



HAL
open science

Limitation strategies for high-order discontinuous Galerkin schemes applied to an Eulerian model of polydisperse sprays

Katia Ait-Ameur, Mohamed Essadki, Marc Massot, Teddy Pichard

► To cite this version:

Katia Ait-Ameur, Mohamed Essadki, Marc Massot, Teddy Pichard. Limitation strategies for high-order discontinuous Galerkin schemes applied to an Eulerian model of polydisperse sprays. *ESAIM: Mathematical Modelling and Numerical Analysis*, 2025, 59 (5), pp.2349-2383. <10.1051/m2an/2025057>. <hal-04374640v2>

HAL Id: hal-04374640

<https://hal.science/hal-04374640v2>

Submitted on 6 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Limitation strategies for high-order discontinuous Galerkin schemes applied to an Eulerian model of polydisperse sprays

K. Ait-Ameur¹, M. Essadki², M. Massot³, and T. Pichard³

¹INRIA, team LEMON / IMAG, Univ. Montpellier, CNRS, 860 Rue Saint Priest, 34095 Montpellier Cedex 5, France

²The MathWorks, 2 Rue de Paris, 92196 Meudon

³Centre de Mathématiques Appliquées, CNRS, École polytechnique, Institut Polytechnique de Paris, Route de Saclay, 91128 Palaiseau Cedex, France

January 6, 2025

Abstract

In this paper, we tackle the modeling and numerical simulation of polydisperse sprays. Starting from a kinetic description for point particles, we focus on an Eulerian high-order geometric method of moment (GeoMOM) in size and consider a system of partial differential equations on a vector of successive fractional size moments of order 0 to $N/2$, $N > 2$, over a compact size interval. These moments correspond to physical quantities, which can be interpreted in terms of the geometry of the interface at small scale. There exists a stumbling block for the usual approaches using high-order moment methods resolved with high-order numerical methods: the transport algorithm does not naturally preserve the moment space. Indeed, reconstruction of moments by polynomials inside computational cells can create N -dimensional vectors which can fail to be moment vectors. We thus propose a new approach, as well as an algorithm, which is high-order in space with limited numerical diffusion, including at the boundaries of the state space, where a specific study is proposed. The main contribution of this work is the design and analysis of a high-order scheme preserving the bounds on the velocity, the moment space and capturing void and δ -shocks solutions. We show that such an approach is competitive compared to second order finite volume schemes, where limiters generate numerical diffusion and clipping at extrema. An accuracy study assesses the order of the method as well as the low level of numerical diffusion on structured meshes. We focus in this paper on cartesian meshes and 2D test cases are presented where the accuracy and efficiency of the approach are assessed.

Introduction

The present work aims at proposing and analyzing a high-order numerical scheme for a system of weakly hyperbolic conservation laws modelling sprays of droplets. In practice, this system is constructed by first considering a collisionless kinetic equation on a distribution function of droplet ([49, 21]), and then by extracting the first few moments with respect to the kinetic variables. The resulting system is underdetermined and it is closed using a quadrature-based approach ([46, 42, 12]). This construction has been widely used (see e.g. the previous work [23, 22, 17] and references therein) to model clouds of spherical droplets. More recently, this approach has also been exploited in [40] for the modelling of multi-scale flow with non-spherical liquid inclusion at small scale.

The considered moment system consists of a pressureless gas dynamics (PGD) system augmented with conservation laws on geometric moments. Therefore the study of the moment system exploit the one of the PGD: First, it is necessary to look for solutions in a weak sense ([6, 4, 7, 50, 5]) because, even with

reasonable smooth initial and boundary conditions, this problem may involve measures, so-called delta-shocks transported with the flow (see also [57, 37, 35]). Second, the uniqueness of the weak solution is only ensured under additional constraints ([6]). Among those constraints, two properties of the initial and boundary values are preserved through space and time: the density remains non-negative and the velocity satisfies a maximum principle, which is closely related to the total variation diminishing (TVD) property. Therefore, these two properties need also to be preserved by numerical schemes for stability reasons. Enforcing the positivity of the density in the PGD corresponds to enforcing that the solution remains in a convex set, called realizability domain or moment set (as it is the set of moments of the non-negative distributions ; [18, 54, 36, 19, 32]), for the moment model. Two types of solutions are difficult to capture by most numerical approaches, those involving concentration of the solution in a spatial point (delta-shock solution), and those involving void regions where the solution is zero as such a value belongs to the boundary of the realizability domain.

At the numerical level, a first order approach for the PGD preserving positivity of the density and the maximum principle on the velocity has been first proposed in [3] based on kinetic interpretation of the PGD, so-called kinetic finite-volume (KFV) scheme. It has been extended to second order in [8] using linear reconstructions with slope limiters to preserve the TVD property on the velocity, and to the considered moment model in [23]. However its construction is restricted to linear reconstructions, therefore to second order schemes usually involving clipping of extrema, since such a reconstruction can be interpreted as a convex combination of the value at each boundary of a cell, while higher order reconstructions do not satisfy such a property. When associated with a strong stability preserving (SSP) Runge Kutta (RK) time discretization ([28, 56]), the discontinuous Galerkin approach (DG ; [14, 15, 20]) has been shown to be a good alternative to construct high-order discretisations for hyperbolic systems with discontinuous solutions.

The DG method produces accurate results if the solution is smooth or contains (relatively) weak discontinuities, otherwise significant oscillations and nonlinear instabilities may occur. To avoid such difficulties with numerical oscillations, the DG method needs to be accompanied by a limitation procedure such as minmod [16], artificial viscosity [41], total variation diminishing [31], weighted essentially nonoscillatory (WENO) [45, 52] techniques or extrema preserving limitations [63, 61]. In addition, there are other techniques for bound-preserving limiters, such as flux corrected transport algorithms [2, 47, 29, 38] and monolithic convex limiting approaches [30, 53]. There have been intensive studies on positivity-preserving and maximum-principle-satisfying methods. The genuinely high-order maximum-principle-satisfying DG method has been proposed in [63, 61] for scalar hyperbolic equations. This procedure has been rapidly developed for different problems ever since, for the Euler equations [64, 65], Navier-Stokes equations [62], shallow water equations [59] and fluid flow in porous media [13], among others. Exploiting a quadrature interpretation of the cell reconstructions, a pointwise limitation has been suggested in this framework. This limitation is rather simple to use, both in term of implementation and to obtain theoretical estimates. However, those estimates are only obtained under the constraint that the solution remains away from the void regions, even if this specific limit is interesting and frequently encountered in applications. Such a scheme has also been tested for the PGD in [60] and for another moment model in [55], which does not involve void or concentrated solutions.

The present work aims at constructing, analyzing and testing some limitation strategies for high-order RKDG discretisations applied to the considered weakly hyperbolic geometric moment system, thus combining the difficulties of moment space preservation within the framework PGD solutions, which can be potentially singular or involve void regions. More specifically, we study the impact of such a limitation in the vicinity of void regions. In practice, the considered limitations can be interpreted as projections of the discrete solution onto the set of admissibility, that is the set of vectors satisfying both the bounds on the velocity and the realizability condition on the geometric moments, while preserving the mean value of the solution in a cell. Two choices of projections towards the cell mean are focused on: one consists in enforcing first the realizability then the overall admissibility, the other consists in projecting all the solution in a single step. These two projections show different behavior in the void region in terms of accuracy and of numerical diffusion.

The paper is organized as follows. In the next section, the construction of the moment model from a kinetic description is recalled and we present its main features in more details. In Section 2, the DG scheme

is presented and the discrete versions of the constraints are specified. Section 3 presents limitation strategies to preserve the realizability and the velocity bounds and their behavior in the vacuum limit is tested on a few test cases. A numerical study is provided in Section 4 to illustrate accuracy and the behavior of the numerical scheme with delta-shocks or vacuum solutions. The last section is devoted to concluding remarks.

1 High-order geometric moment modelling

We present here the construction of the system we aim at solving numerically in the next section, and analyze the properties of its solution we need to preserve at the numerical level.

1.1 Construction of the moment system

The considered system is obtained by evaluating the moments with respect to velocity and size variables of a kinetic model.

1.1.1 Kinetic description

The spray of droplets is described by a Number Density Function (NDF) f , such that $f(t, x, S, v) dx dS dv$ represents the probable number of droplets located in x , with size S and velocity v . The NDF satisfies a Williams-Boltzmann equation ([58]), that models the transport of a spray carried by a gaseous flow

$$\partial_t f + \operatorname{div}_x(vf) = 0. \quad (1)$$

This model is a simple toy problem, but we aim at modelling a more realistic physics. This can be achieved in two manners from (1). First, we can enrich the description of the droplet, typically having a more precise geometry of the droplets, by considering a more complex phase space than only the size variable $S \in \mathbb{R}^+$ (e.g. modelling mean curvatures, temperature or oscillation of the droplets; see e.g. [40, 39, 23, 22]). Second, considering other physical effects such as the drag and evaporation of the droplet (see e.g. [44, 22, 23] and references therein) can simply be modeled through additional terms in (1), potentially depending on these additional variables. However the present contribution in terms of numerical methods naturally extends to these more complex models as the main difficulties arise at the numerical level from the resolution of the transport operator.

1.1.2 Velocity moments

The kinetic phase space is composed of the size variable S and the velocity v . Concerning the velocity variable, we characterize the velocity dependence of f by two quantities: its first two velocity moments

$$\rho = \int_{\mathbb{R}^d} f dv, \quad q = \int_{\mathbb{R}^d} v f dv,$$

that correspond to a density and a momentum. Extracting those moments from (1) yields (formally) the system

$$\begin{aligned} \partial_t \rho + \operatorname{div}_x q &= 0, \\ \partial_t q + \operatorname{div}_x \left(\int_{\mathbb{R}^d} v v^T f dv \right) &= 0, \end{aligned}$$

In order to close this system, and for simplicity, we make the hypothesis that all droplets at a location x are transported at the same velocity u . This corresponds to approximating the distribution f by

$$f(t, x, S, v) \approx \rho(t, x, S) \delta(v - u(t, x)). \quad (2)$$

With such an approximation, we obtain a closed system

$$\partial_t \rho + \operatorname{div}_x(\rho u) = 0, \quad (3a)$$

$$\partial_t q + \operatorname{div}_x(qu^T) = 0, \quad (3b)$$

$$q = \rho u. \quad (3c)$$

This system corresponds to the pressureless gas dynamics (PGD) system, that has been widely studied in the literature (see e.g. [3, 6, 9] and references therein), and where the variable S appears as a parameter. Remark that the absence of pressure in our approach can be justified by the absence of collisions in the underlying kinetic model (1), that is the infinite Knudsen limit, which is valid in many realistic configurations (see e.g. [43]).

1.1.3 Size moments

Concerning the size variable S , we use a fractional moment method ([23]). We use the half order moments $\mathbf{m} = (m_0, m_{1/2}, m_1, m_{3/2})^T$

$$m_\alpha(t, x) = \int_0^1 S^\alpha \rho(t, x, S) dS, \quad (4)$$

where the maximum and minimum admissible sizes are chosen to be $S_{min} = 0$ and $S_{max} = 1$ for simplicity. The reason for this choice of \mathbf{m} is that we retrieve from those moments the following geometric quantities commonly used in the context of separated phase modeling ([21])

$$\Sigma \hat{G} = 4\pi m_0, \quad \Sigma \hat{H} = 2\sqrt{\pi} m_{1/2}, \quad \Sigma = m_1, \quad \alpha = \frac{1}{6\sqrt{\pi}} m_{3/2}, \quad (5)$$

where the α is the volume fraction of liquid, Σ is the interfacial area density and \hat{G} and \hat{H} are respectively the densities of Gauss and mean curvatures averaged over the surface of a droplet ([40, 39, 23, 22]).

Eventually, we extract the moments with respect to $b(S) := (S^0, S^{1/2}, S^1, S^{3/2})^T$ from (3a) and the moment with respect to S only from (3b) to obtain

$$\partial_t U + \operatorname{div}_x(Uu^T) = 0, \quad m_1 u = q, \quad (6)$$

where $U = (\mathbf{m}^T, q^T)^T$ with the moment vector $\mathbf{m} = (m_0, m_{1/2}, m_1, m_{3/2})^T$. The surface area density $\Sigma = m_1$ acts like a density in this model, and the moment against S of (3b) was used to construct q . Following [23], this choice is more relevant when considering drag or evaporation effects, that are not considered in the present work but are part of future projects.

1.2 Properties of the moment system

The considered properties are presented for a 1D version of (6) as such a system can be decomposed into two subsystems widely studied in the literature. The analysis of those 1D subsystems provides some constraints on the solution, our numerical scheme has to satisfy. Eventually, we extend those constraints in a multi-D framework.

1.2.1 Bounds on u

The first subsystem rewritten in 1D is the PGD system on m_1 and $q = (m_1 u)$ which simply consists in a 1D version of (6) where the vector \mathbf{m} is replaced by m_1 , or of (3) replacing ρ by m_1 independent of S .

This system has been analyzed in [3, 8, 6] and we recall a few results here. One specificity of the PGD (3) is the possible appearance of so-called δ -shocks in the solution. It consists of a Dirac measure of mass m_1 transported at velocity u . For this purpose, one needs to focus on solutions in the following weak sense (see [3]).

Definition 1. A couple $(m_1, q) \in C(]0, T[; \mathcal{M}_{loc}(\mathbb{R}))^2$ is a duality solution to the 1D equation (3) if

- The momentum $q \ll m_1$ is absolutely continuous with respect to the mass.
- The mass $m_1 \geq 0$ is non-negative.
- There exists a velocity $u \in L^\infty(]0, T[\times \mathbb{R})$ and a function $a \in L^1_{loc}(]0, T[)$ such that
 - One-sided-Lipschitz condition: $\partial_x u \leq a$.
 - Weak solution: For all $\phi, \psi \in C_c^\infty(]0, T[\times \mathbb{R})$, then

$$\int m_1(\partial_t \phi + u \partial_x \phi) = 0,$$

$$\int q(\partial_t \psi + u \partial_x \psi) = 0.$$

- Representation of u : $m_1 u = q$ a.e. with respect to the measure m_1 .

Remark 1. • In this definition, the velocity $u \equiv \frac{dq}{dm_1}$ corresponds to the Radon-Nikodym derivative of $q \equiv (m_1 u)$ with respect to m_1 on its support $Supp(m_1)$. It is therefore $u(t, \cdot) \in L^\infty(dm_1(t, \cdot))$ for all $t > 0$. Only its definition on $Supp(m_1)$ matters, but this function can be extended in the complement $\mathbb{R} \setminus Supp(m_1)$ into some $L^\infty(]0, T[\times \mathbb{R})$ function (see the notion of universal representative in [3, 6, 4]).

- The requirements on u allow to define a unique characteristic curve X in the sense of Filippov ([50, 26]), i.e. an absolutely continuous function $X(x, t_0) \in W^{1, \infty}(\mathbb{R}^+)$ of t satisfying

$$X(t; x, t_0) = \int_{t_0}^t u(X(\tau; x, t_0)) d\tau, \quad X(t_0; x, t_0) = x.$$

For the numerical application in the next section, we exploit the following property.

Property 1 ([4]). Consider a duality solution (m_1, q) to the 1D equation (3) with initial data m_1^0 and q^0 and denote $u = \frac{dq}{dm_1}$ and $u^0 = \frac{dq^0}{dm_1^0}$. Denote the essential infimum and supremum of a function f with respect to the measure μ on the interval I by $\inf_I^\mu f$ and $\sup_I^\mu f$. Define

$$u_{\inf} := \inf_{\mathbb{R}}^{m_1^0} u^0, \quad u_{\sup} := \sup_{\mathbb{R}}^{m_1^0} u^0.$$

- **Global bound on u :** For all $t > 0$, the velocity u satisfies

$$\begin{aligned} u_{\inf} &\leq \inf_{\mathbb{R}}^{m_1(t, \cdot)} u(t, \cdot), \\ \sup_{\mathbb{R}}^{m_1(t, \cdot)} u(t, \cdot) &\leq u_{\sup}, \end{aligned} \tag{7}$$

- **Local bound on u :** For all $0 < \tau < t$, and all non-empty interval $x \in I^0 = (c, d)$ with $c < d$, define the interval

$$I^\tau = (c - u_{\sup} \tau, d - u_{\inf} \tau).$$

Then the velocity u satisfies

$$\begin{aligned} \inf_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot) &\leq \inf_{I^0}^{m_1(t, \cdot)} u(t, \cdot), \\ \sup_{I^0}^{m_1(t, \cdot)} u(t, \cdot) &\leq \sup_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot). \end{aligned} \tag{8}$$

These constraints are closely related to the total variation diminishing property (TVD; [31]) on velocity, which is also proved to be satisfied by u in [3]. Numerical schemes violating the discrete equivalent of this property can trigger oscillations or overshoots around discontinuities. For this purpose, we rewrite the previous condition.

Property 2. These constraints can be rewritten in terms of the solution (m_1, q) by:

- **Global bound on u :** For all $t > 0$, and all non-empty intervals (c, d) with $c > d$, the solution (m_1, q) satisfies

$$u_{\inf} \int_c^d dm_1(t, \cdot) \leq \int_c^d dq(t, \cdot) \leq u_{\sup} \int_c^d dm_1(t, \cdot). \quad (9)$$

- **Local bound on u :** For all $0 < \tau < t$, and all non-empty intervals $x \in I^0 = (c, d)$ with $c < d$, the solution (m_1, q) satisfies

$$\left(\inf_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot) \right) \int_c^d dm_1(t, \cdot) \leq \int_c^d dq(t, \cdot) \leq \left(\sup_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot) \right) \int_c^d dm_1(t, \cdot). \quad (10)$$

Especially, one observes that the integral of the solution $\left(\int_c^d dm_1(t, \cdot), \int_c^d dq(t, \cdot) \right)$ over an interval (c, d) always belongs to a closed convex cone defined as the intersection of three half spaces: $m_1 \geq 0$ and the two ones defined either globally with u_{\inf} and u_{\sup} in (9) or locally in (10).

Also, the velocity u can not create new local extrema in x and the local minima, resp. maxima, increase, resp. decrease. Following the definition of [31], the velocity is therefore monotonicity preserving and TVD.

1.2.2 Preservation of initial data set

A second subsystem rewritten in 1D yields

$$\partial_t \mathbf{m} + \partial_x (\mathbf{m}u) = 0, \quad (11)$$

where $\mathbf{m} = (m_0, m_{1/2}, m_1, m_{3/2})^T$ is the vector of moments of n with respect to the basis functions $b(S) = (1, S^{1/2}, S, S^{3/2})^T$ and the velocity u is the one found in the previous paragraph from the PGD. We consider again weak solutions in the sense:

Definition 2. A set of functions $\mathbf{m} \in C([0, T[; \mathcal{M}(\mathbb{R}))^4$ is a duality solution to (11) if every of its components m_α , for $\alpha = 0, 1/2, 1, 3/2$, satisfies for all $\phi \in C_c^\infty([0, T[\times \mathbb{R})$

$$\int m_\alpha (\partial_t \phi + u \partial_x \phi) = 0.$$

The considered solutions preserve the initial states:

Proposition 1. Suppose that $m_\alpha^0 \ll m_1^0$ are absolutely continuous with respect to m_1^0 for all $\alpha = 0, 1/2, 1, 3/2$ and that $u(t, \cdot) \in L^\infty(dm_1(t, \cdot))$ is obtained from a duality solution to the 1D equation (3). Then the duality solutions to (11) satisfy for all borel set B and $t > 0$,

$$\mathbf{m}(t, B) \in \text{Cone}(\{\mathbf{m}^0(y), y \in \mathbb{R}\}),$$

where $\text{Cone}(\cdot)$ is the convex cone pointed at the origin generated by all initial \mathbf{m}^0 .

Proof. This results from the method of characteristics in the sense of Filippov ([26, 50]). Remark that all the components m_α follow the same characteristic curve which provides the result. The requirement that m_α^0 is dominated by m_1^0 simply provides that $\text{Supp}(m_\alpha^0) \subset \text{Supp}(m_1^0)$ and therefore following the characteristics provides $\text{Supp}(m_\alpha(t, \cdot)) \subset \text{Supp}(m_1(t, \cdot))$ and ensures the uniqueness of the velocity in the support $\text{Supp}(m_\alpha(t, \cdot))$. \square

1.2.3 Moment set and Hankel determinants

This initial set is encompassed into a larger set, that is the set of moments of n with respect to $b(S)$, also called the realizability domain. This set of moments is often studied when constructing moment models because an important part of the physics is put into the nonlinear source terms, which are only defined and numerically evaluated under realizability constraints. This constraint extends the constraint $m_1 \geq 0$ in Definition 1.

Definition 3. The set of all realizable moment vectors or realizability domain yields

$$\mathcal{R} = \left\{ \int_0^1 b(S) d\mu(S), \quad \mu \in \mathcal{M}([0, 1]) \right\}.$$

Property 3 ([1, 18, 48]). This set is a closed convex cone characterized by its extremal points

$$\mathcal{R} = \text{Cone}(\{b(S), \quad S \in [0, 1]\}).$$

This set is characterized by numerical constraints following Hausdorff problem.

Proposition 2. The vector $\mathbf{m} = (m_0, m_{1/2}, m_1, m_{3/2})^T \in \mathbb{R}^4$ is realizable if the following matrices are symmetric non-negative

$$H^1 = \begin{pmatrix} m_{1/2} & m_1 \\ m_1 & m_{3/2} \end{pmatrix}, \quad H^2 = \begin{pmatrix} m_0 - m_{1/2} & m_{1/2} - m_1 \\ m_{1/2} - m_1 & m_1 - m_{3/2} \end{pmatrix}, \quad (12a)$$

and all their components are non-negative. This is equivalent to requiring their trace and determinants are non-negative, which reformulates $h_i(U) \geq 0$ with

$$\begin{aligned} h_1(U) &= m_{1/2} + m_{3/2}, & h_2(U) &= m_0 - m_{1/2} + m_1 - m_{3/2}, \\ h_3(U) &= m_{1/2}m_{3/2} - m_1^2, & h_4(U) &= (m_0 - m_{1/2})(m_1 - m_{3/2}) - (m_{1/2} - m_1)^2. \end{aligned} \quad (12b)$$

Proof. The fractional realizability condition simply follows from the solution of Hausdorff moment problem ([18, 36, 54, 19, 32, 48]) after using a change of variable $S = r^2$ in the integration. \square

1.2.4 Admissible set

Eventually, we call globally, resp. locally, admissible the set of function that satisfy both the global, resp. local, bounds on u of Section 1.2.1 and the realizability of Section 1.2.3.

Definition 4. Suppose that the solution U to (6) for all time $0 < \tau < t$. The set of globally admissible vectors is

$$\mathcal{A}_{t,[a,b]}^{glob,\tau} := \{(\mathbf{m}^T, q)^T \text{ such that } \mathbf{m} \in \mathcal{R}, \quad m_1 u_{inf} \leq q \leq m_1 u_{sup}\},$$

and the set of locally admissible vectors

$$\mathcal{A}_{t,[a,b]}^\tau := \{(\mathbf{m}^T, q)^T \text{ such that } \mathbf{m} \in \mathcal{R}, \quad m_1 u_{inf}^\tau \leq q \leq m_1 u_{sup}^\tau\}, \quad (13a)$$

$$u_{inf}^\tau(t) = \inf_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot), \quad u_{sup}^\tau(t) = \sup_{I(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot), \quad (13b)$$

Again, these sets are defined as intersection of closed convex cones and are therefore closed convex cones.

These sets are defined such that the vectors obtained by integrating the solution $U(t, \cdot)$ to (6) over any interval I belong to this set, i.e.

$$\left(\int_c^d dU(t, \cdot) \right) \in \mathcal{A}_{t,[a,b]}^\tau \subset \mathcal{A}_{t,[a,b]}^{glob,\tau}$$

and where the velocity u in the Definitions 4 is the one resulting from the solution U over $0 < \tau < t$.

1.2.5 Extension to multi-D problems

We extend the framework for multi-D problems, but the analysis is left for future work. When considering a problem of spatial dimension $d > 1$, we consider duality solutions under the following sense.

Definition 5. A couple $U = (\mathbf{m}^T, q^T)^T \in C([0, T[; \mathcal{M}_{loc}(\mathbb{R}^d))^{4+d}$ is a duality solutions of (6) if:

- Every component $q_i, m_\alpha \ll m_1$ for all $i = 1, \dots, d$ and $\alpha = 0, 1/2, 1, 3/2$, is absolutely continuous w.r.t. m_1 .
- The vector $\mathbf{m} \in \mathcal{R}$ is realizable m_1 -a.e.
- There exists $u \in L^\infty([0, T[\times \mathbb{R}^d)^d$ such that
 - Weak solution: For all component U_i and $\forall \phi \in C_c^\infty([0, T[\times \mathbb{R}^d)$, then

$$\int U_i (\partial_t \phi + u^T \nabla_x \phi) = 0.$$

- Representation of u : $m_1 u = q$ is satisfied m_1 -a.e.

Remark that an entropy condition à la Oleinik has been present in Definition 1 in the 1D case. It is missing in the present extension, and we do not perform a proper analysis of the multi-D model. A first result in this direction has been proposed in [7] for the transport equation and extension of this work to (6) is left for future work.

Assuming that there still exists a unique Filippov characteristics passing at every $(t, x) \in]0, T[\times \mathbb{R}^d$, then we extend in multi-D:

- **Globally directional velocity bound:** The bound (9) on the velocity applies in every direction, i.e. for all $n \in \mathbb{S}^d$, all non-empty set $C \subset \mathbb{R}^d$, for all $t > 0$,

$$(u^T n)_{\inf} \int_C dm_1(t, \cdot) \leq \int_C d(n^T q)(t, \cdot) \leq (u^T n)_{\sup} \int_C dm_1(t, \cdot),$$

where $(u^T n)_{\inf} = \inf_{\mathbb{R}^d}^{m_1^0} (u^T n)$ and $(u^T n)_{\sup} = \sup_{\mathbb{R}^d}^{m_1^0} (u^T n)$.

- **Locally directional velocity bound:** The bound (10) extends for all $n \in \mathbb{S}^d$, all non-empty set $C^0 \subset \mathbb{R}^d$, for all $t > 0$,

$$\left(\inf_{C(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot) \right) \int_C dm_1(t, \cdot) \leq \int_C dq(t, \cdot) \leq \left(\sup_{C(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot) \right) \int_C dm_1(t, \cdot),$$

where

$$C^\tau := \{x - u, \quad x \in C, \quad (u^T n) \in [\tau(u^T n)_{\sup}, \tau(u^T n)_{\inf}]\}.$$

- The realizability property $\mathbf{m} \in \mathcal{R}$ naturally extends if this vector is transported along characteristic curves.

Eventually the admissible set (13) extends into

$$\mathcal{A}_{t,C}^\tau := \{(\mathbf{m}^T, q)^T \text{ such that } \mathbf{m} \in \mathcal{R}, \quad m_1(u^T n)_{\inf}^\tau \leq q^T n \leq m_1(u^T n)_{\sup}^\tau\}, \quad (14a)$$

$$(u^T n)_{\inf}^\tau(t) = \inf_{C(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot), \quad (u^T n)_{\sup}^\tau(t) = \sup_{C(t-\tau)}^{m_1(t-\tau, \cdot)} u(t-\tau, \cdot). \quad (14b)$$

1.2.6 Discussion on numerical difficulties

Two types of difficulties are focused on in the numerical section below:

- **The appearance of void:** at certain locations, the moment solution can become zero. In this limit, one of the inequality in (12) becomes an equality, and therefore the moment $\mathbf{m} = (m_0, m_{1/2}, m_1, m_{3/2})^T \in \partial\mathcal{R}$ belongs to the boundary of its admissible set. Furthermore, the velocity u is ill-defined in this limit because it only appears multiplied by m_1 in (6). Such an issue has been illustrated in [3] through a 1D test case for the PGD, that we extend in the present framework into

$$m_\alpha^0(x) = \int_0^1 S^\alpha dS = (1 + \alpha)^{-1}, \quad q^0(x) = m_1^0(x) \times \begin{cases} 0.5 & \text{for } x > 0, \\ -0.5 & \text{for } x < 0. \end{cases} \quad (15)$$

This m_α^0 is in the interior of \mathcal{R} and generates a void region in finite time as the solution simply yields

$$m_\alpha(t, x) = (1 + \alpha)^{-1} \times (1 - \mathbf{1}_{[-0.5t, 0.5t]}(x)), \quad q(t, x) = 0.5 \times \begin{cases} 0.5 & \text{for } x > 0.5t, \\ -0.5 & \text{for } x < -0.5t. \end{cases}$$

- **The appearance of δ -shocks:** the solution may contain Dirac measures that are propagated with the flow. Again such a solution has been exhibited in [3] through a test case rewritten into

$$m_\alpha^0(x) = (1 + \alpha)^{-1}, \quad q^0(x) = m_1^0(x) \times \begin{cases} -0.5 & x > 0, \\ 0.5 & x < 0. \end{cases} \quad (16)$$

It generates a δ -shock as the solution yields

$$m_\alpha(t, x) = (1 + \alpha)^{-1}(1 + 2t\delta_0(x)), \quad q(t, x) = q^0(x).$$

The following section presents a high-order scheme preserving the bounds on the velocity, the moment space and capturing void and δ -shocks solutions.

2 RKDG scheme preserving admissibility

Following the derivation of the model (6) from the kinetic model (1), seen as a Galerkin approximation with respect to the kinetic variable, we extend here this construction with a strong stability preserving (SSP) Runge-Kutta discontinuous Galerkin (RKDG) scheme. A special focus is given to the imposition of the admissibility enforcement on the numerical solution.

2.1 Discontinuous Galerkin (DG) space discretization

In order to construct the space discretization of (6), remark first that the moment method of Section 1 is already a Galerkin approximation of (1) with respect to the kinetic variables (S, v) using \mathbf{b} for the test functions and (2) for approximation function. We extend the Galerkin approximation with space discretization.

For simplicity, we use a Cartesian grid $D = \bigcup_e \Omega_e$ where Ω_e is a product of intervals of the form $[x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}]$. Define polynomial test functions $g \in \mathbb{P}_r(\Omega_e)$ of degree r with respect to the space variable x and compute the integral

$$\begin{aligned} 0 &= \int_{\mathbb{R}^d} \int_0^1 \int_{\Omega_e} g(x) \mathbf{b}(S, v) (\partial_t f + \text{div}_x(vf))(t, x, v, S) dx dS dv \\ &= \frac{d}{dt} \int_{\mathbb{R}^d} \int_0^1 \int_{\Omega_e} g(x) \mathbf{b}(S, v) f(t, x, v, S) dx dS dv \\ &\quad - \int_{\mathbb{R}^d} \int_0^1 \int_{\Omega_e} (\nabla_x g(x)^T v) \mathbf{b}(S, v) f(t, x, v, S) dx dS dv \\ &\quad + \int_{\mathbb{R}^d} \int_0^1 \int_{\partial\Omega_e} g(x) (v^T n(x)) \mathbf{b}(S, v) f(t, x, v, S) dx dS dv, \end{aligned} \quad (17)$$

where n denotes the outgoing normal to the boundary $\partial\Omega_e$. In the spirit of [3, 8], following the characteristic curves suggests to decompose f in the last integral into two parts coming from both side of the interface. Considering $x \in \Gamma_{ee'} = \Omega_e \cap \Omega_{e'}$ on the interface between Ω_e and $\Omega_{e'}$ and denoting n the normal directed toward $\Omega_{e'}$, this corresponds to writing

$$(v^T n(x))f(t, x, v, S) = \lim_{\substack{y \rightarrow x \\ y \in \Omega_e}} (v^T n(y))_+ f(t, y, v, S) + \lim_{\substack{y \rightarrow x \\ y \in \Omega_{e'}}} (v^T n(y))_- f(t, y, v, S), \quad (18)$$

where $a_{\pm} = (a \pm |a|)/2$ designate the positive or negative part of a . Now, following (2), we approximate f by f_h defined as

$$f_h(t, x, v, S) = \rho_h(t, x, S) \prod_j \delta_{u_{j,h}(t,x)}(v_j) = \sum_e \rho^e(t, x, S) \prod_j \delta_{u_j^e(t,x)}(v_j) \mathbf{1}_{\Omega_e}(x), \quad (19)$$

where the subscript h refers to the functions defined by parts and the superscript e refers to the functions in the cell Ω_e . The functions u^e and ρ^e are chosen such that $x \mapsto U^e(t, x) = ((\mathbf{m}^e)^T, (q^e)^T)(t, x)^T \in \mathbb{P}_r(\Omega_e)^{4+d}$ are polynomial function of space in every cell Ω_e , where

$$U^e(t, x) = \int_0^1 \int_{\mathbb{R}^d} \mathbf{b}(S, v) \rho^e(t, x, S) \prod_j \delta_{u_j^e(t,x)}(v_j) dS = \int_0^1 \mathbf{b}(S, u^e(t, x)) \rho^e(t, x, S) dS \quad (20)$$

are polynomials of degree r over the spatial cell Ω_e . Eventually, only the moment equation is solved, and the fact that $U_e(x)$ is polynomial is sufficient for the construction¹.

Now, we choose the g_j such that it forms a basis of polynomials, which is orthogonal with respect to the L^2 scalar product on Ω_e . Denote x_k some quadrature points. In practice, we simply choose in 1D the Gauss-Lobatto points such that they include the boundary of each interval, and they maximize the accuracy in the sense that the space integrals are exact up to degree $2r - 1$. Finally, denote l_k the Lagrange polynomials associated to these quadrature points. This structure (quadrature points and Lagrange polynomials) is simply tensorized in multi-D (see [61, 64] and references therein). Reinjecting it in (17) provides

$$0 = M \frac{d}{dt} U - F(U) + E(U), \quad (21a)$$

where the unknown $(U_{j,k})_{j=1,\dots,4+d} = ((\mathbf{m}_k)^T, (q_k)^T)^T$ approximates $U(x_k) = ((\mathbf{m}^e)^T, (q^e)^T)(x_k)^T$ at the quadrature points x_k and

$$\left(M \frac{dU}{dt} \right)_{i,j} = \sum_k \left(\int_{\Omega_e} g_i(x) l_k(x) dx \right) \frac{dU_{j,k}}{dt}, \quad (21b)$$

$$F(U)_{i,j} = \sum_k \left(\int_{\Omega_e} g_i(x) \nabla_x l_k(x)^T dx \right) u_k U_{j,k}, \quad (21c)$$

where the velocity u_k satisfies $m_{1,k} u_k = q_k \approx q(x_k) = m_1(x_k) u(x_k)$. It is a scalar in 1D, or a vector in multi-D of the same size as $\nabla_x l_k(x)$. The case $m_{1,k} = 0$, which corresponds to the zero mass $m_1 = 0$ case, will be treated in the next section.

For the exchange term E , the boundary $\partial\Omega_e = \cup_{e'} \Gamma_{ee'}$ of the cell is splitted into the interfaces $\Gamma_{ee'} = \Omega_e \cap \Omega_{e'}$ and one remarks that the quadrature points x_k along $\Gamma_{ee'}$ in Ω_e are identical to those on the other side, along $\Gamma_{ee'}$ in $\Omega_{e'}$. Therefore, reinjecting (18) in the last integral of (17) and using the approximation (19) leads to (see also [3, 8])

$$E(U)_{i,j} = \sum_{\substack{e' \text{ s.t.} \\ \Gamma_{ee'} \neq \emptyset}} \sum_{\substack{k \text{ s.t.} \\ x_k \in \Gamma_{ee'}}} \left(\int_{\Gamma_{ee'}} g_i(x) l_k(x) dx \right) \times \left((u_k^T n(x_k))_+ U_{j,k} + (u_k^T n(x_k))_- U_{j,k'} \right), \quad (21d)$$

¹Such an approximation can be achieved in several ways, as such underlying ρ^e and u^e exist. For instance, $\rho^e(t, x, S) = \sum_{i=1}^4 \alpha_i(t, x) \delta_{S_i}(S)$, $u_j^e(t, x) = p_j(t, x) / (\sum_{i=1}^4 \alpha_i(t, x) S_i)$, with $x \mapsto \alpha_i(t, x), p_j(t, x) \in \mathbb{P}_r(\Omega_e)$ and some distinct fixed $S_i \in (0, 1)$ provides such U_e . Other choices are possible and this formula does not impact our construction.

where the index k refers to the quadrature points along the edge $\Gamma_{ee'}$ in Ω_e and the index $k' \neq k$ corresponds to the quadrature point in $\Omega_{e'}$ at the same location $x_{k'} = x_k \in \Gamma_{ee'}$. In 1D, this exchange terms reduces to the kinetic fluxes [3].

2.2 Numerical admissibility constraint

The conditions (12) and (8) satisfied by the duality solution at the continuous level need to be transposed at the discrete level.

Realizability: For the realizability condition (12), despite providing a discretized version of the underlying kinetic model, the transposition of this condition at the discrete level is essentially driven by the applications we have in mind. Indeed, violating a discrete version of the pointwise realizability criteria (12) would not affect the precision nor the stability of the scheme, as a non-realizable vector would simply be transported at velocity u . The main motivation for imposing this constraint arise from the additional physical effects discussed in Section 1.1.1 that would require a strong imposition of this property.

Velocity bounds: The preservation of monotonicity in the solution, and therefore the bounds on velocity, is closely related to the total variation diminishing property (TVD; [31]). In 1D, the velocity u has been shown to be total variation diminishing (TVD) in [3] at the continuous level and preserving the bounds on the total variation at the numerical level is essential for stability. These bounds are applied at the cell level Ω_e to the cell mean values.

Formulation of the requirement at the cell level: First, denote \bar{U}_e the integral of the approximation U^e in the cell Ω_e . Using the appropriate Gauss-Lobatto quadrature weights $\omega_k > 0$

$$\bar{U}_e = \frac{1}{|\Omega_e|} \int_{\Omega_e} U^e(x) dx = \sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k U_k \approx \frac{1}{|\Omega_e|} \int_{\Omega_e} \int_0^1 \int_{\mathbb{R}^d} \mathbf{b}(v, S) f(x, v, S) dv dS dx,$$

which approximates the moments of f in the spatial cell Ω_e with $|\Omega_e|$, the element size. Assuming that $U(t^n, \cdot)$ is of the form (20) at time t^n with positive $\rho(t^n, \cdot)$, then we expect the integral of exact solution to satisfy $\bar{U}_e^{n+1} \in \mathcal{A}_{t^n, \Omega_e}^{\Delta t}$ as defined in (14) for all cell Ω_e . This rewrites:

$$\begin{aligned} \bar{\mathbf{m}}_e^{n+1} &\in \mathcal{R}, & \bar{q}_{j,e}^{n+1} &\in [\bar{m}_{1,e}^{n+1} (u^T e_j)_{\min,e}^n, \bar{m}_{1,e}^{n+1} (u^T e_j)_{\max,e}^n], \\ (u^T n)_{\min,e}^n &= \min_{\substack{e' \text{ s.t.} \\ \Omega_e \cap \Omega_{e'} \neq \emptyset}} (\bar{u}_{e'}^n)^T n, & (u^T n)_{\max,e}^n &= \max_{\substack{e' \text{ s.t.} \\ \Omega_e \cap \Omega_{e'} \neq \emptyset}} (\bar{u}_{e'}^n)^T n, \end{aligned}$$

where \bar{u}_e^n satisfies $\bar{m}_{1,e}^n \bar{u}_e^n = \bar{q}_e^n$. Assuming that the time step satisfies a condition of the form $\Delta t \leq \max_e \bar{u}_e^n R(\Omega_e)$ where the radius $R(\Omega_e)$ is the maximum distance between two points of Ω_e , this corresponds to imposing (12) and (8) to the approximate solution f^h averaged in a cell Ω_e at time t^{n+1} .

Formulation of the requirement at the node level: Exploiting the positivity of the quadrature weights of Gauss-Lobatto and the convexity of the admissible set, $\bar{U}_e^{n+1} \in \mathcal{A}_{t^n, \Omega_e}^{\Delta t}$ holds if $U_k^{n+1} \in \mathcal{A}_{t^n, \Omega_e}^{\Delta t}$ holds for every quadrature point $x_k \in \Omega_e$, or equivalently

$$\mathbf{m}_k^{n+1} \in \mathcal{R}, \quad q_{j,k}^{n+1} \in [m_{1,k}^{n+1} (u^T e_j)_{\min,e}^n, m_{1,k}^{n+1} (u^T e_j)_{\max,e}^n], \quad (23)$$

implies (22). For the applications below, we impose (23) and we rewrite the discrete admissible set as

$$\mathcal{A}_e^n := \{U \text{ such that } \mathbf{m} \in \mathcal{R}, \quad m_1 (u^T e_j)_{\min,e}^n \leq q_j \leq m_1 (u^T e_j)_{\max,e}^n\}, \quad \mathcal{A}^n := \prod_e \mathcal{A}_e^n, \quad (24)$$

using velocities $(u^T e_j)_{\min,e}^n$ and $(u^T e_j)_{\max,e}^n$ given at a time step t^n in the cell Ω_e . Especially, at every time step, the discrete admissible set is a closed convex cone defined using the solution U^{n-1} at the previous time step.

Imposition of the requirement: A numerical scheme preserves admissibility if

$$U^n \in \mathcal{A}^{n-1} \quad \Rightarrow \quad U^{n+1} \in \mathcal{A}^n.$$

At each time step, we enforce realizability (12b) of the moments vector (it boils down to the positivity of the density for the PGD system [3]) and local maximum principle for the velocity (8). Global bound preserving limiter can limit the approximation order reached by the numerical scheme (see [66]). These constraints (13) holds true at the continuous level and we seek to satisfy them at the discrete level (23). If admissibility is lost at some time step and at some quadrature point during the simulation, we correct this numerical solution in the following way.

For $U \notin \mathcal{A}^n$, we define a corrected value ($\mathcal{P}^n U$) as a projection onto \mathcal{A}^n and it needs to satisfy:

- $(\mathcal{P}^n U) \in \mathcal{A}^n$ is admissible,
- For all Ω_e ,

$$\sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k U_k = \bar{U}_e = \sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k (\mathcal{P}^n U)_k. \quad (25)$$

The second criteria aims at imposing conservativity of the scheme. Indeed, the numerical scheme satisfied by the cell-averaged quantities \bar{U}_e^{n+1} with the time discretizations of the next subsection can be rewritten in a conservative (finite volume) manner. Correcting the numerical solution at every time step such that (25) holds, this numerical scheme can still be written in a conservative way, but with modified fluxes.

In practice, we also expect the correction $\|U - \mathcal{P}^n U\|$ to be as small as possible in order not to deteriorate the accuracy of the scheme. Especially, the restriction $\mathcal{P}^n|_{\mathcal{A}^n}$ to the admissible set must be the identity. This additional error is studied in the next two sections.

Following the work of [61, 63, 64, 65], we consider corrections of the form:

$$(\mathcal{P}^n U)_{i,k} = \theta_{i,e}^n U_{i,k} + (1 - \theta_{i,e}^n) \bar{U}_{i,e} \quad \text{for all } k \text{ such that } x_k \in \Omega_e, \quad (26)$$

where the convex combination parameters $\theta_{i,e}^n \in [0, 1]$ are such that $(\mathcal{P}^n U) \in \mathcal{A}^n$. They can be different for the different components i of the vector $U = (\mathbf{m}^T, q^T)^T$, for different cells Ω_e and for different time t^n . But they are the same for every quadrature points $x_k \in \Omega_e$ among the same cell in order to satisfy the conservativity property (25).

Property 4. Corrections of the form (26) preserve the conservativity property (25) of the scheme (21).

Proof. Let $(\mathcal{P}^n U)_{i,e}$ be the modified polynomial after corrections in cell Ω_e , for the different components i , such that

$$(\mathcal{P}^n U)_{i,e}(x_k) = (\mathcal{P}^n U)_{i,k}, \quad x_k \in \Omega_e.$$

Then, the cell average $\overline{(\mathcal{P}^n U)_{i,e}}$ satisfies:

$$\begin{aligned} \overline{(\mathcal{P}^n U)_{i,e}} &= \sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k (\mathcal{P}^n U)_{i,k} = \sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k (\theta_{i,e}^n U_{i,k} + (1 - \theta_{i,e}^n) \bar{U}_{i,e}) \\ &= \theta_{i,e}^n \left(\sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k U_{i,k} \right) + (1 - \theta_{i,e}^n) \bar{U}_{i,e} \left(\sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k \right) = \bar{U}_{i,e}, \end{aligned}$$

since $\theta_{i,e}^n$ does not depend on the quadrature point x_k and the quadrature weights verifies $\sum_{\substack{k \text{ s.t.} \\ x_k \in \Omega_e}} \omega_k = 1$. \square

2.3 SSP Runge-Kutta time discretization

Concerning the time discretization, we exploit the strong stability preserving (SSP) Runge-Kutta (RK) framework ([28, 56]). Such schemes have been originally designed to preserve the TVD property while going higher order.

2.3.1 Admissibility of the explicit Euler-finite volume scheme

Rewriting the DG semi-discretization (21a) of the last subsection under the form

$$\frac{dU}{dt} = \mathcal{G}(U) = M^{-1}(F - E)(U), \quad (27)$$

then the explicit Euler iteration reads

$$\tilde{U}^{n+1} = U^n + \Delta t \mathcal{G}(U^n), \quad (28)$$

that can be corrected if $\tilde{U}^{n+1} \notin \mathcal{A}^n$. The explicit Euler time discretizations combined with first order kinetic finite volume spatial discretization reads (28) where \mathcal{G} is defined in (27) with M , F and E defined in (21) and using only one quadrature point per cell and therefore only one function $l_1(x) = 1 = g_1(x)$ per cell (such that $\nabla_x l_1 = 0$). This yields especially $U_k = \bar{U}_e$ for the unique $x_k \in \Omega_e$. Then decomposing the term $\bar{U}_e = U_k^n$ over the dual mesh (in Cartesian geometry) as in discrete duality finite volume (DDFV ; see e.g. [10, 11]), the scheme (28) reduces to (for k such that $x_k \in \Omega_e$)

$$\begin{aligned} U_k^{n+1} &= U_k^n - \Delta t \sum_{\substack{x_{k'} \in \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \frac{((u_k^n)^T n_{kk'})_+ U_k^n + ((u_{k'}^n)^T n_{kk'})_- U_{k'}^n}{\Delta x_{kk'}} \\ &= U_k^n - \Delta t \sum_{\substack{x_{k'} \in \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \frac{\mathcal{F}_{kk'}^n - \mathcal{F}_{k'k}^n}{\Delta x_{kk'}}, \quad \mathcal{F}_{kk'}^n = ((u_k^n)^T n_{kk'})_+ U_k^n, \end{aligned} \quad (29)$$

and, remarking that $a_+ = -(-a)_-$, we used $u_k^n \equiv q_k^n / (m_1)_k^n$ and

$$\Delta x_{kk'} = \Delta x_{k'k} = \|x_k - x_{k'}\|, \quad n_{kk'} = -n_{k'k} = \frac{x_k - x_{k'}}{\Delta x_{kk'}},$$

First, we rewrite the result (similar to those in [3, 8]) in the simple first order framework:

Proposition 3 (Admissibility preservation of explicit Euler-kinetic finite volume scheme). *Suppose that $U^n \in \mathcal{A}^{n-1}$ and that*

$$\frac{\Delta t \max_k \|u_k^n\|}{\min_{\substack{(x_k, x_{k'}) \in \Omega_e \times \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \Delta x_{k'k}} \leq 1. \quad (30)$$

Then the update $U^{n+1} \in \mathcal{A}^n$ computed with (29) on a Cartesian grid is admissible.

Proof. From (29), the scheme rewrites (for k such that $x_k \in \Omega_e$)

$$U_k^{n+1} = \left(1 - \Delta t \sum_{\substack{x_{k'} \in \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \frac{((u_k^n)^T n_{kk'})_+}{\Delta x_{kk'}} \right) U_k^n + \Delta t \sum_{\substack{x_{k'} \in \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \frac{((u_{k'}^n)^T n_{k'k})_+}{\Delta x_{k'k}} U_{k'}^n,$$

and, under the CFL condition (30) and since the mesh is Cartesian, it is a convex combination of the previous values U_k^n . Therefore, it preserves any convex set, and especially it preserves the admissibility property. \square

Remark 2. This construction extends out of the Cartesian framework under a modified CFL condition. Cartesian grids were found sufficient for the present applications.

2.3.2 Admissibility of the cell-averaged value of the Euler-DG scheme

A first step toward the high-order RKDG scheme is the explicit Euler-DG scheme, namely high order in space and first order in time. It is given in (28). When using a higher order DG spatial discretization, i.e. with a large set of polynomials $(l_i)_i$ and $(g_i)_i$, then the scheme a priori does not preserve admissibility. However, for the admissibility of the full SSPRK-DG scheme in the next paragraph, we only need the admissibility of its cell-averaged value, i.e. if $\bar{U}_e^n \in \mathcal{A}_e^{n-1}$ for all e , then we need $\bar{U}_e^{n+1} \in \mathcal{A}_e^n$ for all e , where

$$\bar{U}_e^{n+1} = \sum_{x_k \in \Omega_e} w_k \tilde{U}_k^{n+1}, \quad (31)$$

and \tilde{U}_k^{n+1} is given in (28).

Proposition 4 (Admissibility preservation of the scheme on the cell-averaged of the SSPRK-DG scheme). *Suppose that $\bar{U}^n \in \mathcal{A}^{n-1}$ and that*

$$\frac{\Delta t \max_k \|u_k^n\|}{\left(\min_{\substack{(x_k, x_{k'}) \in \Omega_e \times \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \Delta x_{kk'} \right) \times \min_k w_k} \leq 1. \quad (32)$$

Then the updated value $\bar{U}^{n+1} \in \mathcal{A}^n$ of the cell-averaged U^{n+1} computed with (31) on a Cartesian grid is admissible.

Proof. We simply reformulate proofs from the literature (e.g. in Theorem 2.1 [64], Theorem 4.1 [60], Theorem 12.33 [33]) to the present framework. The scheme (31) can be written as a combination of Euler-KFV steps

$$\bar{U}^{n+1} = \sum_{x_k \in \Omega_e} w_k \left(U_k^n - \frac{\Delta t}{w_k} \sum_{\substack{x_{k'} \in \tilde{\Omega}_{e'} \\ \text{s.t. } \tilde{\Gamma}_{ee'} \neq \emptyset}} \frac{((u_k^n)^T n_{kk'})_+ U_k^n + ((u_{k'}^n)^T n_{k'})_- U_{k'}^n}{\Delta x_{kk'}} \right), \quad (33)$$

by redecomposing the mesh into another one with cells $\tilde{\Omega}_e$ and edges $\tilde{\Gamma}_{ee'}$ constructed such that there is only one quadrature point x_k per cell at its center. This scheme was shown to preserve the admissibility from one step to another in Proposition 3 under the considered CFL condition, by dilatating Δt by a coefficient $1/w_k$. \square

This implies that the cell-averaged values of the numerical solution remain admissible at every iteration, but the value at each quadrature point might leave the admissible set anyway. This is circumvented by the projection in the next section.

2.3.3 Admissibility of the limited SSPRK-DG scheme

When using a higher order DG spatial discretization, i.e. with a large set of polynomials $(l_i)_i$ and $(g_i)_i$, then the scheme a priori does not preserve admissibility. For this reason, the scheme is corrected into

$$U^{n+1} = \mathcal{P}^n \tilde{U}^{n+1}, \quad (34)$$

where the projections $\mathcal{P}^n \tilde{U}^{n+1} \equiv \mathcal{P}^n(\tilde{U}^{n+1}, \bar{U}^{n+1})$ toward the cell-averaged values \bar{U}^{n+1} are defined in the next section. For this projection to exist this cell averaged value needs to be admissible $\bar{U}^{n+1} \in \mathcal{A}$, and the scheme satisfied by these values needs to preserve admissibility, as described in the last paragraph.

Eventually, we construct a corrected m -stage SSPRK-DG scheme under the form

$$\begin{cases} U^{(0)} &= U^n, \\ U^{(j)} &= \sum_{i=0}^{j-1} \alpha_{j,i} \mathcal{P}^n \left(U^{(i)} + \Delta t \frac{\beta_{j,i}}{\alpha_{j,i}} \mathcal{G}(U^{(i)}), \overline{U^{(i)} + \Delta t \frac{\beta_{j,i}}{\alpha_{j,i}} \mathcal{G}(U^{(i)})} \right) \quad \text{for } j = 1, \dots, m, \\ U^{n+1} &= U^{(m)}, \end{cases} \quad (35)$$

where the coefficients $\alpha_{j,i}, \beta_{j,i}$ are chosen such that the RK scheme is strong stability preserving (SSP ; [28]), i.e. they satisfy

$$\alpha_{j,i}, \beta_{j,i} \geq 0, \quad \sum_{i=0}^{j-1} \alpha_{j,i} = 1. \quad (36)$$

The cell averaged scheme $\overline{U^{(i)} + \Delta t \frac{\beta_{j,i}}{\alpha_{j,i}} \mathcal{G}(U^{(i)})}$ in (35) corresponds exactly to (33) by replacing Δt by $\Delta t \frac{\beta_{j,i}}{\alpha_{j,i}}$.

In this work, the RK method is always chosen of the same order as the one in space. These schemes are chosen to be SSP in order to preserve admissibility (see Proposition 5 below), and up to fourth order to use only standart explicit SSPRK methods. Higher order discretizations would require further development beyond the present objectives. The parameters α and β used in Section 4 below are given in Propositions 3.1, 3.2 and 3.3 in [27] for the second, third and fourth order SSPRK schemes. These coefficients are chosen such that if no correction is needed, that is if $\mathcal{P}^n = Id$ for all j in (35), then $(\bar{U}_e^{n+1} - \bar{U}_e^n)/\Delta t = O(\Delta t^p)$ at a certain order p .

Proposition 5 (Admissibility preservation of the SSPRK-DG scheme). *Suppose that $U^n \in \mathcal{A}^{n-1}$ and that*

$$\frac{\Delta t \max_{i,j} \frac{|\beta_{j,i}|}{\alpha_{j,i}} \times \max_k \|u_k^n\|}{\left(\min_{\substack{(x_k, x_{k'}) \in \Omega_e \times \Omega_{e'} \\ \text{s.t. } \Gamma_{ee'} \neq \emptyset}} \Delta x_{k'k} \right) \times \min_k w_k} \leq 1. \quad (37)$$

Then the update $U^{n+1} \in \mathcal{A}^n$ computed with (35), with parameters α, β satisfying (36), with a projection \mathcal{P}^n onto the admissible set \mathcal{A}^n , on a Cartesian grid, is admissible.

Proof. This scheme is again constructed as a convex combination of admissible values $\mathcal{P}^n(U, \bar{U}) \in \mathcal{A}^n$, well-defined under the given CFL condition according to Proposition 4. \square

3 Projection methods

Two projections of the form (26) are constructed and studied in this section in order to enforce admissibility. They are defined locally as functions \mathcal{P} depending on $U \in \mathbb{R}^{4+d}$ and $\bar{U} \in \mathcal{A}_e^n$, corresponding respectively to the quadrature value U_k^n and the cell value \bar{U}_e^n . Rewrite (26) generically

$$\mathcal{P}(U, \bar{U})_i = \theta_i U_i + (1 - \theta_i) \bar{U}_i, \quad (38)$$

where the coefficients θ_i are defined below.

3.1 Projections enforcing realizability and TVD of the velocity

We define values of θ_i to enforce the admissibility requirement. This condition rewrites $h_i(U) \geq 0$ for $i = 1, \dots, 4 + 2d$, the first four correspond to realizability, the last to the TVD of u . In practice, we impose $h_i(U) \geq \varepsilon > 0$ with a small value of ε in order to avoid admissibility loss due to round-off error. This parameter is fixed to 10^{-12} such that the admissibility domain is not too reduced, this was found sufficient for our applications.

3.1.1 Realizability projection

First, we define the projection to enforce realizability:

Definition 6 (Realizability projection). For $\bar{\mathbf{m}}$ and \mathbf{m} such that $h_i(\bar{U}) > \varepsilon$ and $h_i(U) < \varepsilon$ for some $i = 1, \dots, 4$ (defined in (12b)), the realizability projection is defined by

$$\mathcal{P}^{real}(U, \bar{U}) = \theta_i U + (1 - \theta_i) \bar{U},$$

such that $\theta_i \in [0, 1]$ satisfy $h_i(\mathcal{P}U) = \varepsilon$ and thus are given by:

$$\theta_1 = \frac{\varepsilon - h_1(\bar{U})}{h_1(U - \bar{U})}, \quad \theta_2 = \frac{\varepsilon - h_2(\bar{U})}{h_2(U - \bar{U})}, \quad (39a)$$

$$\theta_3 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad \theta_4 = \frac{-\tilde{b} + \sqrt{\tilde{b}^2 - 4\tilde{a}\tilde{c}}}{2\tilde{a}} \quad (39b)$$

$$\text{with } \begin{cases} a = h_3(U - \bar{U}), & c = h_3(\bar{U}) - \varepsilon, \\ b = \bar{m}_{1/2}(m_{3/2} - \bar{m}_{3/2}) + (m_{1/2} - \bar{m}_{1/2})\bar{m}_{3/2} - 2\bar{m}_1(m_1 - \bar{m}_1), \\ \tilde{a} = h_4(U - \bar{U}), & \tilde{c} = h_4(\bar{U}) - \varepsilon, \\ \tilde{b} = (\bar{m}_0 - \bar{m}_{1/2})(m_1 - \bar{m}_1 - m_{3/2} + \bar{m}_{3/2}) \\ \quad + (\bar{m}_1 - \bar{m}_{3/2})(m_0 - \bar{m}_0 - m_{1/2} + \bar{m}_{1/2}) \\ \quad - 2(\bar{m}_{1/2} - \bar{m}_1)(m_{1/2} - \bar{m}_{1/2} - m_1 + \bar{m}_1). \end{cases}$$

One verifies that each $\theta_i \in [0, 1]$ is well-defined as long as $h_i(U) < \varepsilon$ and $h_i(\bar{U}) > \varepsilon$.

3.1.2 TVD projection

Similarly, the constraints (24) on q for the TVD property on u rewrites $h_i(U) \geq 0$ for $i = 5, \dots, 4 + 2d$ with:

$$h_{4+2j-1}(U) = q_j - m_1(u^T e_j)_{min}, \quad h_{4+2j}(U) = m_1(u^T e_j)_{max} - q_j. \quad (40)$$

Definition 7 (TVD projection). For \bar{U} and \hat{U} such that $h_i(\bar{U}), h_i(\hat{U}) > \varepsilon$ for $i = 1, \dots, 4$ and $h_j(\bar{U}) > \varepsilon$ and $h_j(\hat{U}) < \varepsilon$ for $j = 5, \dots, 4 + 2d$, we define the TVD projection of the form

$$\mathcal{P}^{TVD}(\hat{U}, \bar{U}) = \theta_j \hat{U} + (1 - \theta_j) \bar{U},$$

where $\theta_j \in [0, 1]$ is such that $h_j(\mathcal{P}\hat{U}) = \varepsilon$. This yields

$$\theta_j = \frac{\varepsilon - h_j(\bar{U})}{h_j(\hat{U} - \bar{U})}. \quad (41)$$

Remark that this definition requires $\hat{\mathbf{m}} \in \mathcal{R}$ to be realizable for the coefficients $\theta_j \in [0, 1]$. Therefore, in order to construct projections over the admissibility domain, we need to combine these two projections \mathcal{P}^{real} and \mathcal{P}^{TVD} , and the realizability one needs always to be imposed first. Two combinations of those projections are considered in the next paragraph.

3.2 Assembling the projections

3.2.1 Minimal projection

We first introduce a projection \mathcal{P}_{min} that is only used for accuracy study.

Definition 8 (Minimal projection). The closest admissible vector to a non-admissible one $U \notin \mathcal{A}_e^n$ (red curve on Fig. 1):

$$\mathcal{P}_{min}(U, \bar{U}) = \operatorname{argmin}_{V \in \mathcal{A}_e^n} \|U - V\| \quad (42)$$

The projection \mathcal{P}_{\min} does not depend on \bar{U} . Especially, it can a priori not be written in the form (38) in order to preserve the conservativity Property 4. For this reason, it is only exploited below for accuracy studies, but it is not used for the limitation in the DG scheme.

By definition, this projection $\mathcal{P}_{\min}(U, \bar{U}) \in \mathcal{A}_\varepsilon^n$ is the admissible vector the closest to $U \notin \mathcal{A}_\varepsilon^n$. Especially, for the accuracy study below, we exploit the estimation between the projection $\mathcal{P}_{\min}(U, \bar{U})$ and the exact solution U^{ex}

$$\|\mathcal{P}_{\min}(U, \bar{U}) - U^{ex}\| \leq \|\mathcal{P}_{\min}(U, \bar{U}) - U\| + \|U - U^{ex}\| \leq 2\|U - U^{ex}\|. \quad (43)$$

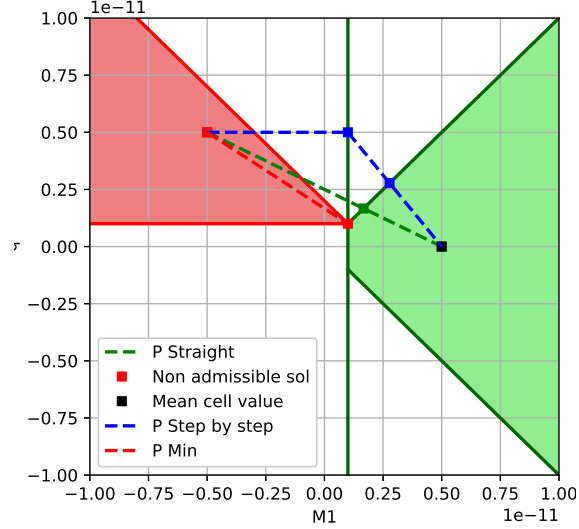


Figure 1: Cut in the (m_1, q) plane of the admissible set and the minimal, step-by-step and straight projections of a non-admissible vector with $\varepsilon = 10^{-12}$.

3.2.2 Step-by-Step projection

Definition 9 (Step-by-step projection). Inspired of [61, 64, 65, 63, 66, 62, 59], the projection is performed in two steps. We first project \mathbf{m} onto $\hat{\mathbf{m}} \in \mathcal{R}$, then project $\hat{U} = (\hat{\mathbf{m}}^T, q^T)^T$ onto $\mathcal{A}_\varepsilon^n$.

- Define first $\hat{U} = \mathcal{P}_1(U, \bar{U})$ such that (blue curve in Fig. 1)

$$\hat{U} = \mathcal{P}_1(U, \bar{U}) = (\theta^1 \mathbf{m}^T + (1 - \theta^1) \bar{\mathbf{m}}^T, q^T)^T, \quad \text{where } \theta^1 = \min_{i=1, \dots, 4} \theta_i, \quad (44a)$$

where $\theta_i = 1$ if $h_i(U) \geq \varepsilon$ or equals (39) otherwise, such that $\mathcal{P}^1(U, \bar{U})$ satisfies (12). Remark that \mathcal{P}_1 has its first components \mathbf{m} equal to those of \mathcal{P}^{real} , but has different q .

- Then, project $\hat{U} = \mathcal{P}_1(U, \bar{U})$ onto $\mathcal{A}_\varepsilon^n$ using the functions h_i for $i = 5, \dots, 4 + 2d$:

$$\mathcal{P}_{SbS}(U, \bar{U}) = \theta^2 \hat{U} + (1 - \theta^2) \bar{U} \quad \text{where } \theta^2 = \min_{i=5, \dots, 4+2d} \theta_i, \quad (44b)$$

where $\theta_i = 1$ if $h_i(\hat{U}) \geq \varepsilon$ or equals (41) otherwise, such that $\mathcal{P}_{SbS}(U, \bar{U})$ satisfies (23).

Property 5. The second projection (44b) does not alter the realizability property, and the step-by-step projection is a projection onto the admissible set, i.e. if $\bar{U} \in \mathcal{A}_\varepsilon^n$, then $\mathcal{P}_{SbS}(U, \bar{U}) \in \mathcal{A}_\varepsilon^n$ for all U .

Proof. First, one remarks that \mathcal{A}_e^n is a convex cone, since it is the set of moments of positive distributions. This implies that for all i

$$\begin{aligned} h_i(V) \geq 0 &\Rightarrow h_i(\alpha V) \geq 0 \quad \forall \alpha \geq 0, \\ h_i(V), h_i(W) \geq 0 &\Rightarrow h_i(\theta V + (1 - \theta)W) \geq 0 \quad \forall \theta \in [0, 1]. \end{aligned}$$

Especially, since both $h_i(\hat{U}) \geq \varepsilon$ and $h_i(U) \geq \varepsilon$ for $i = 1, \dots, 4$, then the second projection (44b) does not alter the realizability of \mathbf{m} . \square

Equivalently, the projection \mathcal{P}_{SbS} corresponds to fixing $\theta_i = \theta^2$ onto the component q (i.e. $i \geq 5$) and $\theta_i = \theta^1 \theta^2$ (i.e. $i \leq 4$) onto the components \mathbf{m} in (38).

3.2.3 Straight projection

Definition 10 (Straight projection). This projection consists in using the same parameter θ for every component (see green curve in Figure 1)

$$\mathcal{P}_{Str}(U, \bar{U}) = \theta^3 \hat{U} + (1 - \theta^3) \bar{U}, \quad \hat{U} = \theta^1 U + (1 - \theta^1) \bar{U}, \quad (45)$$

where θ^1 is defined in (44a) and θ^3 is such that $\mathcal{P}_{Str}(\hat{U}, \bar{U})$ satisfies (23). It yields $\theta^3 = \min_{i=5, \dots, 4+2d} \theta_i$ with either $\theta_i = 1$ if $h_i(\hat{U}) \geq 0$ or it equals (41) otherwise.

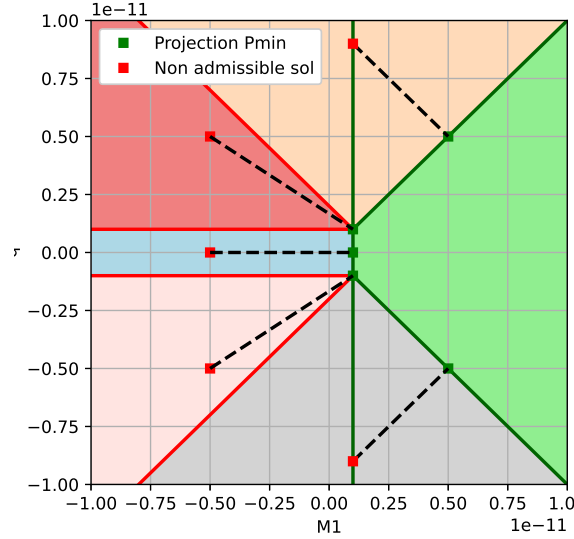


Figure 2: Projection \mathcal{P}_{\min} applied to (m_1, q)

Property 6. These projections are comparable in several regimes:

- When $m_1 \geq \varepsilon$, then $\mathcal{P}_{SbS}(U, \bar{U}) = \mathcal{P}_{Str}(U, \bar{U})$.
- When $0 \leq q \leq \varepsilon u_{max}$ (in blue on Fig. 2), then $\mathcal{P}_{SbS}(U, \bar{U}) = \mathcal{P}_{\min}(U, \bar{U})$.
- When $q \geq \varepsilon u_{max}$ and $m_1 \leq \varepsilon$, then $\mathcal{P}_{SbS}(U, \bar{U}) \neq \mathcal{P}_{Str}(U, \bar{U})$ a priori.

Eventually, both \mathcal{P}_{SbS} and \mathcal{P}_{Str} take the form (38) and can therefore be used in a conservative manner in the DG scheme, while \mathcal{P}_{\min} can not. Their accuracy are studied in Section 3.3.

3.3 Behavior near the vacuum regime

Following the work of [61, 64, 65, 63, 66, 62, 59], one easily derives error estimations with the different limitations under the condition that \bar{U} is far from the boundary $\partial\mathcal{A}_\varepsilon^n$. However, as illustrated in Section 1.2.6, the near vacuum regime $m_1 \rightarrow 0^+$, that corresponds to a case along the boundary $\partial\mathcal{A}_\varepsilon^n$ can be trigger by relatively common configurations and is therefore relevant in many applications. For this reason, we study the properties of $\mathcal{P}_{SbS}(U, \bar{U})$ and $\mathcal{P}_{Str}(U, \bar{U})$ in the limit $\bar{m}_1 \rightarrow 0^+$. For simplicity, we conduct this study in 1D in the (m_1, q) plane where the admissibility property simplifies into:

$$\varepsilon \leq m_1, \quad m_1 u_{\min} + \varepsilon \leq q \leq m_1 u_{\max} - \varepsilon, \quad (46)$$

but this extends to the case where the condition $m_1 \geq \varepsilon$ is replaced by $\mathbf{m} \in \mathcal{R}$.

3.3.1 Mass comparison

First, one observes that the mass is higher with Step-by-Step projection in this limit.

Proposition 6. Consider $U = (m_1, q)^T \notin \mathcal{A}_\varepsilon^n$ and $\bar{U} = (\bar{m}_1, \bar{q})^T \in \mathcal{A}_\varepsilon^n$ such that

$$q < \varepsilon u_{\max} - (\varepsilon - m_1) \frac{\bar{q} - \varepsilon u_{\max}}{\bar{m}_1 - \varepsilon} \quad \text{or} \quad q \geq \bar{q}.$$

Then, writing $\mathcal{P}_{SbS}(U, \bar{U}) = (m_1^{SbS}, q^{SbS})^T$ and $\mathcal{P}_{Str}(U, \bar{U}) = (m_1^{Str}, q^{Str})^T$

$$m_1^{Str} \leq m_1^{SbS}.$$

Proof. Suppose that $m_1 \geq \varepsilon$, then $m_1^{SbS} = m_1^{Str}$.

Suppose that $m_1 < \varepsilon$. One verifies that the first case $q < \varepsilon u_{\max} - (\varepsilon - m_1) \frac{\bar{q} - \varepsilon u_{\max}}{\bar{m}_1 - \varepsilon}$ also provides $m_1^{SbS} = \varepsilon = m_1^{Str}$.

Finally, for the case $q \geq \bar{q}$, comparing the values (44a) and (45) provides the inequality. \square

Remark 3. The vacuum limit corresponds to having $\bar{q} \in [\bar{m}_1 u_{\min}, \bar{m}_1 u_{\max}]$ in the limit $\bar{m}_1 \rightarrow 0^+$. Near this limit, and assuming that the numerical error of the DG scheme is much bigger than ε , the case $q < \bar{q}$ is less probable. Indeed, the size of the interval $[\bar{m}_1 u_{\min}, \bar{m}_1 u_{\max}]$ in which falls \bar{q} is much smaller than the DG error.

This comparison is supported by some numerical examples below and may result in larger numerical diffusion effects on m_1 with \mathcal{P}_{SbS} than with \mathcal{P}_{Str} .

3.3.2 Error comparisons

For a quantitative study of the error, we obtain the following estimations:

Lemma 1. Consider vectors $U \notin \mathcal{A}_\varepsilon^n$, $U^{ex}, \bar{U} \in \mathcal{A}_\varepsilon^n$, and a projection $\mathcal{P}(U, \bar{U}) \in \mathcal{A}_\varepsilon^n$ onto the admissibility domain, using the average value \bar{U} . Then

$$\|\mathcal{P}(U, \bar{U}) - U^{ex}\| \leq \|(\mathcal{P} - \mathcal{P}_{min})(U, \bar{U})\| + 2\|U - U^{ex}\|. \quad (47)$$

Proof. Following the computations of [61, 64, 65, 63, 66, 62, 59],

$$\|\mathcal{P}(U, \bar{U}) - U^{ex}\| \leq \|(\mathcal{P} - \mathcal{P}_{min})(U, \bar{U})\| + \|\mathcal{P}_{min}(U, \bar{U}) - U^{ex}\|,$$

and using (43) provides the result. \square

The discrete solution $U \notin \mathcal{A}_\varepsilon^n$ is supposed to approximate the exact one $U^{ex} \in \mathcal{A}_\varepsilon^n$ with a certain accuracy (typically $\mathcal{O}(\Delta x^p)$ in the DG framework). Then, we need to study the distance $(\mathcal{P} - \mathcal{P}_{min})$ for the two projections \mathcal{P}_{SbS} and \mathcal{P}_{Str} when $\bar{U} \rightarrow (\varepsilon, \varepsilon u)$ with $u \in [u_{\min}, u_{\max}]$.

Concerning \mathcal{P}_{min} , straightforward computations provide:

Lemma 2 (Values of \mathcal{P}_{\min}). • Suppose that (in yellow on Fig. 2)

$$m_1 u_{\max} - q \leq 0 \quad \text{and} \quad m_1 + q u_{\max} \geq \varepsilon(1 + u_{\max}^2),$$

then $\mathcal{P}_{\min}(U, \bar{U})$ is the orthogonal projection onto the axis $q = m_1 u_{\max}$, i.e.

$$\mathcal{P}_{\min}(U, \bar{U}) = (m_1^{\min}, m_1^{\min} u_{\max})^T, \quad m_1^{\min} = \frac{m_1 + q u_{\max}}{1 + u_{\max}}.$$

• Suppose that (in red on Fig. 2)

$$m_1 + q u_{\max} \leq \varepsilon(1 + u_{\max}^2) \quad \text{and} \quad q \geq u_{\max} \varepsilon,$$

then

$$\mathcal{P}_{\min}(U, \bar{U}) = (\varepsilon, \varepsilon u_{\max}).$$

• Suppose that (in blue on Fig. 2)

$$0 \leq q \leq u_{\max} \varepsilon,$$

then $\mathcal{P}_{\min}(U, \bar{U}) = (\varepsilon, q)$.

The ones for negative q can be deduced by symmetry.

Eventually, straightforward computations yield the distance $\|(\mathcal{P} - \mathcal{P}_{\min})(U, \bar{U})\|$ in (47) in the different cases:

Proposition 7 (Error estimation with \mathcal{P}_{SbS}). Assuming that (46) holds, then we have $\hat{U} = (\varepsilon, q)^T$ as defined in (44a) and

$$\mathcal{P}_{SbS}(\hat{U}, \bar{U}) = \theta_q \hat{U} + (1 - \theta_q) \bar{U}, \quad \theta_q = \frac{\bar{m}_1 u_{\max} - \bar{q}}{q - \bar{q} - u_{\max}(\varepsilon - \bar{m}_1)}.$$

Then, suppose additionnally that

• Either (red on Fig. 1)

$$m_1 + q u_{\max} \leq \varepsilon(1 + u_{\max}^2), \quad \text{and} \quad q \geq u_{\max} \varepsilon,$$

then

$$\|(\mathcal{P}_{SbS} - \mathcal{P}_{\min})(U, \bar{U})\| \leq \|(1, u_{\max})\| |1 - \theta_q| (\bar{m}_1 - \varepsilon).$$

Epecially, $\lim_{\bar{U} \rightarrow (\varepsilon, \varepsilon u)} \|(\mathcal{P}_{SbS} - \mathcal{P}_{\min})(U, \bar{U})\| = 0$.

• Or (yellow in Fig. 1)

$$m_1 + q u_{\max} \geq \varepsilon(1 + u_{\max}^2) \quad \text{and} \quad m_1 \leq \varepsilon,$$

then

$$\|(\mathcal{P}_{SbS} - \mathcal{P}_{\min})(U, \bar{U})\| \leq \|(1, u_{\max})\| |\theta_q m_1 + \bar{m}_1(1 - \theta_q) - m_1^{\min}|.$$

Epecially, $\lim_{\bar{U} \rightarrow (\varepsilon, \varepsilon u)} \|(\mathcal{P}_{SbS} - \mathcal{P}_{\min})(U, \bar{U})\| = \|(1, u_{\max})\| |m_1^{\min} - \varepsilon| \neq 0$.

In the last case (yellow in Fig. 1), one needs further assumptions on U to control this error. This is typically done in the DG framework by exploiting the value of an exact solution U^{ex} assumed to be close to U .

Proposition 8 (Error estimation with \mathcal{P}_{Str}). Assuming that (46) holds, then

$$\mathcal{P}_{Str}(U, \bar{U}) = \theta U + (1 - \theta) \bar{U}, \quad \theta = \min(\theta_{m_1}, \theta_q), \quad \theta_{m_1} = \frac{\varepsilon - \bar{m}_1}{m_1 - \bar{m}_1}, \quad \theta_q = \frac{\bar{m}_1 u_{\max} - \bar{q}}{q - \bar{q} - u_{\max}(\varepsilon - \bar{m}_1)}. \quad (48)$$

We distinguish the two cases $\theta = \theta_{m_1}$ and $\theta = \theta_q$ (see Fig. 3) depending on the location of \bar{U} and U . Suppose additionnally that

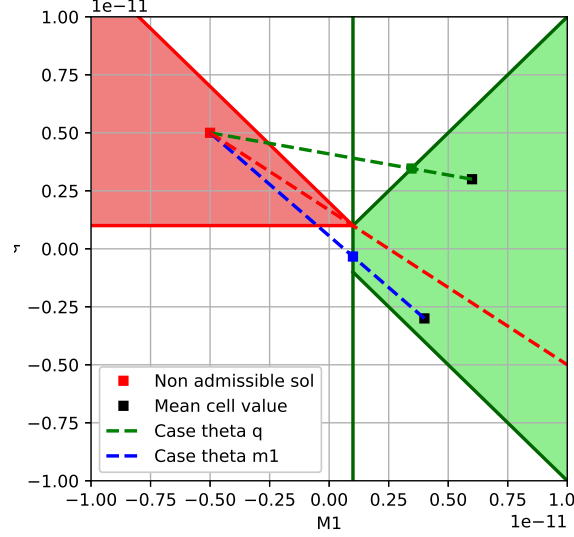


Figure 3: Projection \mathcal{P}_{Str} depending on the location of \bar{U}

- **Case 1:** $\theta = \theta_{m_1}$ and (red in Fig. 1)

$$m_1 + qu_{max} \leq \varepsilon(1 + u_{max}^2), \quad \text{and} \quad q \geq u_{max}\varepsilon,$$

then

$$\|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\| = \varepsilon|u - u_{max}|, \quad u = \frac{\theta q + (1 - \theta)\bar{q}}{\theta m_1 + (1 - \theta)\bar{m}_1}.$$

Epecially, $\|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\|$ is always considered negligible.

- **Case 2:** $\theta = \theta_q$ and (red in Fig. 1)

$$m_1 + qu_{max} \leq \varepsilon(1 + u_{max}^2) \quad \text{and} \quad q \geq u_{max}\varepsilon,$$

then

$$\|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\| \leq \|(1, u_{max})\| |m_1^{Str} - \varepsilon|, \quad m_1^{Str} = \theta m_1 + (1 - \theta)\bar{m}_1.$$

Epecially, using $\varepsilon \leq m_1^{Str} \leq \bar{m}_1$, then $\lim_{\bar{m}_1 \rightarrow \varepsilon} \|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\| = 0$.

- **Case 3:** $\theta = \theta_{m_1}$ and (yellow in Fig. 1)

$$m_1 + qu_{max} \geq \varepsilon(1 + u_{max}^2) \quad \text{and} \quad m_1 \leq \varepsilon,$$

then

$$\|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\| \leq \varepsilon|u - u_{max}| + u_{max}|\varepsilon - m_1^{min}|.$$

Epecially, the first term $\varepsilon|u - u_{max}|$ is negligible, but the second term $|\varepsilon - m_1^{min}|$ is a priori not controlled and one needs again further assumptions on $\mathcal{P}_{min}(U, \bar{U})$.

- **Case 4:** $\theta = \theta_q$ and (yellow in Fig. 1)

$$m_1 + qu_{max} \geq \varepsilon(1 + u_{max}^2) \quad \text{and} \quad m_1 \leq \varepsilon,$$

then

$$\|(\mathcal{P}_{Str} - \mathcal{P}_{min})(U, \bar{U})\| \leq \|(1, u_{max})\| |m_1^{Str} - m_1^{min}|.$$

- **Case 5:** $\theta = \theta_{m_1}$ and (yellow in Fig. 1) $0 \leq q \leq \varepsilon u_{max}$. Then $\mathcal{P}_{Str}(U, \bar{U}) = \mathcal{P}_{min}(U, \bar{U})$.
- **Case 6:** $\theta = \theta_q$ and (yellow in Fig. 1) $0 \leq q \leq \varepsilon u_{max}$. Then

$$\|(\mathcal{P}_{min} - \mathcal{P}_{Str})(U, \bar{U})\| \leq \varepsilon |u - u_{max}| + u_{max} |m_1^{Str} - \varepsilon|.$$

Especially, the first term $\varepsilon |u - u_{max}|$ is negligible and, since $\varepsilon \leq m_1^{Str} \leq \bar{m}_1$, the second term satisfies $\lim_{\bar{m}_1 \rightarrow \varepsilon} |\varepsilon - m_1^{Str}| = 0$.

4 Numerical experiments

The present approach is constructed to reach high-order accuracy (restricted to second order in [8, 34, 25, 24]) while remaining robust with void regime and singularities. This is illustrated in this section through five representative test-cases, with increasing difficulties. The first is a smooth 1D case to study the accuracy of the method. The second and third cases test the robustness of the method respectively when void or δ -shocks appear. The last test cases extend the void and δ -shock studies in a 2D framework. Our numerical results are compared with a kinetic finite volume scheme (KFV; [8, 34, 25, 24]). It is a finite volume scheme using the kinetic fluxes (21d) at the interfaces, together with a MUSCL linear reconstruction in each cell. A minmod limiter based on u and the so-called canonical moments formulation is used. It consists in a non-linear transformation of the realizability domain \mathcal{R} onto $[0, 1]^4$. However, this strategy is limited to second order. The time step satisfies the following CFL condition for the second order kinetic finite volume scheme:

$$\frac{\Delta t}{\Delta x} = \text{CFL} \max_e \left(|\bar{u}_e| + |Du_e| \frac{\Delta x}{2} \right), \quad (49)$$

where $\text{CFL} \leq 1$ is defined by the user. The quantity \bar{u}_e is different from the cell average quantity u_e and is determined depending on the slope Du_e and conservation properties. The slope Du_e is calculated in order to satisfy realizability and maximum principle conditions. The expressions of \bar{u}_e and Du_e are given in the appendix D of [24]. The time step for the RKDG schemes satisfies the CFL condition (37) where the first Gauss-Lobatto weight ω_1 is 1 for the second order RKDG scheme, $\frac{1}{2}$ for the third order scheme and $\frac{5}{9}$ for the fourth order scheme.

4.1 Accuracy study for a 1D spray

The first initial condition yields

$$m_\alpha(x, 0) = \int_0^1 S^\alpha G(x, 1/2) dS = (\alpha + 1)^{-1} G(x, 1/2), \quad q(x, 0) = -m_1(x, 0), \quad (50)$$

for $\alpha = 0, 1/2, 1, 3/2$, where $G(x, x_c) = \exp(-(x - x_c)^2/\sigma^2)$ with $\sigma = 0.1$. Periodic conditions are used at the boundary of the $[0, 1]$ -domain. With such a field, all the moments are transported at velocity $u = -1$. Fig. 4 (left) shows the numerical solution m_0 (the other moments show identical features) with 100 cells at time $t = 2$, i.e. after two periods. Fig. 4 (right) shows the relative l^1 -distance of the numerical solutions m_0^N obtained with the different schemes to the exact solution at final time $t^N = 2$ as a function of Δx

$$\varepsilon(\Delta x) = \frac{\sum_e \sum_q \omega_q |m_{0,q}^N - m_0(x_q, t = 2)|}{\sum_e \sum_q \omega_q m_0(x_q, t = 2)}.$$

Both limitations of Section 3 have no impact for this simulation and give identical result. Therefore, only one is shown and the RKDG schemes yield the desired order. The KFV scheme [34] with the minmod limitation yields a slightly slower convergence than the theoretical second order, and the maximum is clipped. Without limitation, the RKDG and the KFV schemes yield a similar underlying spatial reconstruction. The difference in the behavior is due to the considered limitation and the minmod limitation on the canonical moments affects more the local extrema of the solution.

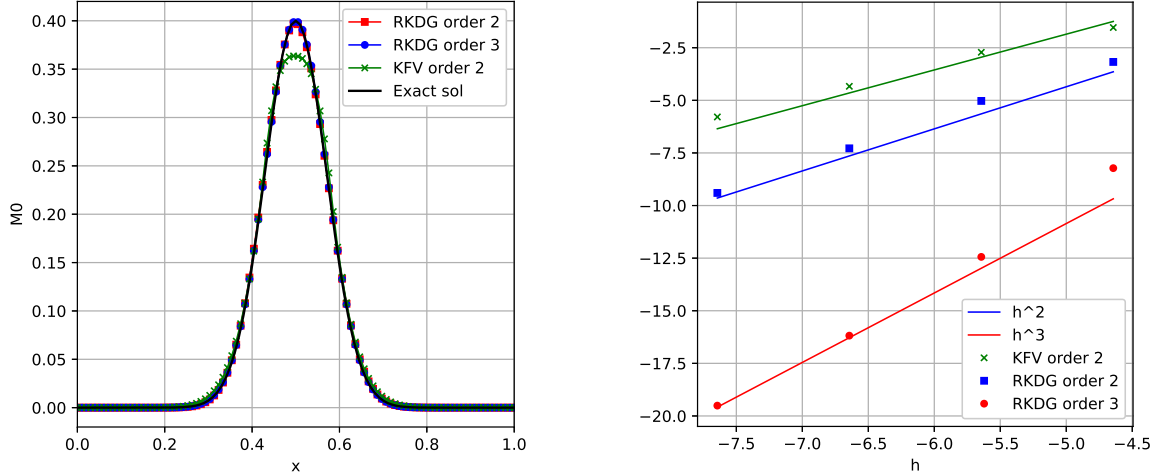


Figure 4: Moment m_0 (left) obtained with the 2nd and 3rd order RKDG schemes and the 2nd order KJV scheme [34] at $t^N = 2$ with the initial condition (50). Relative l^1 distances (right) between the numerical solutions m_0^N at $t^N = 2$ and the exact solution.

4.2 Vacuum test case

We extend a test case from [8] in the present moment framework (see (15)). We use the initial condition derived in (15):

$$m_\alpha(x, 0) = (\alpha + 1)^{-1}, \quad q(x, 0) = m_1(x, 0) \times \begin{cases} -0.4 & \text{if } x < 0.5 \text{ or } x > 1.8, \\ 0.4 & \text{if } 0.5 < x < 1, \\ 1.4 - x & \text{if } 1 < x < 1.8. \end{cases} \quad (51)$$

The first gap in the initial velocity is meant to trigger the low density region and periodic boundary conditions are used.

Fig. 5 displays the numerical solutions m_0 (left) and q (right) at $t^N = 0.5$ with 100 cells. Vacuum is created at the location of the velocity discontinuity in the initial condition. The KJV scheme is more diffusive in this region. As in the previous test case, the profile in the higher density region is sharper with the third order scheme and more diffused with the KJV scheme. However, the RKDG schemes present larger overshoots on the sides of this profile and in the middle of this region. This was already observed in [8]. This effect reduces when raising the order of accuracy. Other slope limiters allow to damp these oscillations near shocks like WENO-based limiters [67]. This limitation procedure is achieved by using a troubled-cell indicator [51] to identify the cells that require reconstruction. Then a polynomial reconstruction is made with a WENO-inspired approach (see [67]). In our work, we seek to ensure realizability and local maximum principle and we identify the cells violating these two constraints to correct them with the slope limiters (26).

A further study is carried out with Fig. 6 and Table 1. Here, we seek to compare the behavior near vacuum of the RKDG scheme when the step-by-step (44a)-(44b) or the straight projection (45) is applied in the limitation. Fig. 6 show zooms on this void region with the second and fourth order RKDG schemes where the two projections are compared. Table 1 presents the minimal value of m_1 obtained with the second, third and fourth order RKDG schemes where the two projections are also compared. The safety parameter used in the limitations (39) and (41) is set to $\varepsilon = 10^{-12}$. At second order with coarse mesh size, the two projections show identical minimal value of m_1 , this corresponds to the case where only one constraint is violated i.e. the realizability and positivity of m_1 . In all the other simulations, the straight limitation \mathcal{P}^{Str} (45) shows lower values of m_1 compared to the step-by-step one \mathcal{P}^{SbS} (44b), as observed in Subsection 3.3.1. The even orders, i.e. odd order polynomial reconstructions, show lower values of m_1 in the void region with the

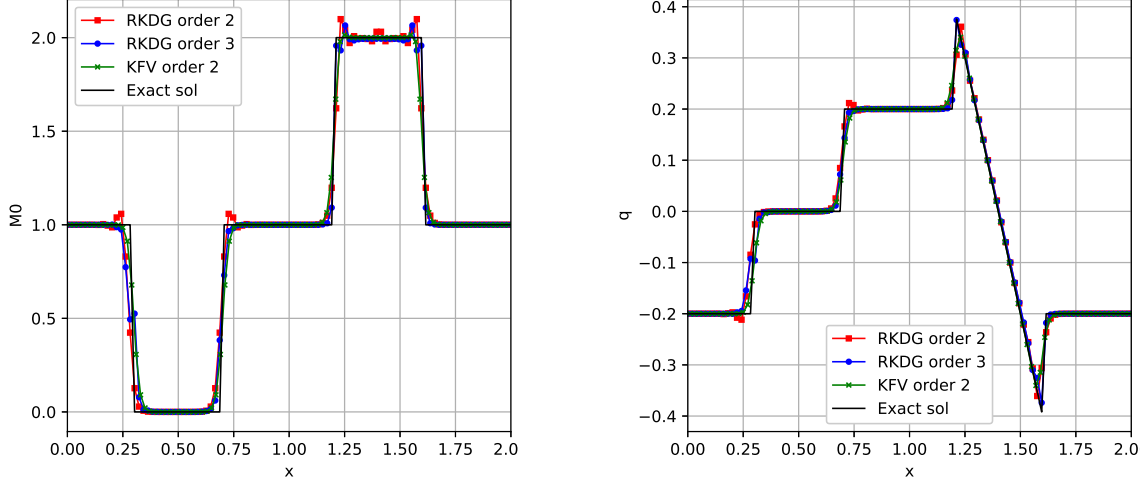


Figure 5: Moments m_0 (left) and q (middle) obtained with the RKDG schemes and the KJV scheme at $t = 0.4$ with the initial conditions (51).

| Mesh size | RKDG Order 2 | | RKDG Order 3 | |
|-----------|-------------------------|-------------------------|------------------------|------------------------|
| | \mathcal{P}_{Sbs} | \mathcal{P}_{Str} | \mathcal{P}_{Sbs} | \mathcal{P}_{Str} |
| $N = 25$ | 8.015×10^{-3} | 8.015×10^{-3} | 1.225×10^{-2} | 1.225×10^{-2} |
| $N = 50$ | 6.497×10^{-5} | 6.497×10^{-5} | 6.297×10^{-5} | 6.510×10^{-5} |
| $N = 100$ | 4.221×10^{-9} | 4.221×10^{-9} | 2.813×10^{-6} | 8.459×10^{-6} |
| $N = 200$ | 1.775×10^{-12} | 1.403×10^{-12} | 6.881×10^{-7} | 5.688×10^{-8} |

| Mesh size | RKDG Order 4 | | KJV Order2 |
|-----------|-------------------------|--------------------------|--------------------------|
| | \mathcal{P}_{Sbs} | \mathcal{P}_{Str} | |
| $N = 25$ | 3.188×10^{-3} | 5.096×10^{-7} | 5.859×10^{-3} |
| $N = 50$ | 2.281×10^{-3} | 1.448×10^{-5} | 6.103×10^{-5} |
| $N = 100$ | 2.460×10^{-7} | $10^{-12} = \varepsilon$ | 5.587×10^{-9} |
| $N = 200$ | 2.060×10^{-12} | $10^{-12} = \varepsilon$ | $10^{-12} = \varepsilon$ |

Table 1: Minimal value of m_1 with the different schemes and limitations in the vacuum test case.

straight projection than with the step-by-step one. As a summary, the present high-order scheme remains robust when considering solutions with very low m_1 .

4.3 1D δ -shock test case

We extend another test case from [8] using:

$$\begin{aligned}
 m_\alpha(x, 0) &= \int_0^1 S^\alpha \left(G(x, x_1) + \frac{4}{3} G(x, x_2) \mathbf{1}_{[0.5, 1]}(S) \right) dS \\
 &= (\alpha + 1)^{-1} \left(G(x, x_1) + \frac{4}{3} (1 - 0.5^{\alpha+1}) G(x, x_2) \right), \tag{52a}
 \end{aligned}$$

$$q(x, 0) = m_1(x, 0) (\mathbf{1}_{\mathbb{R}^-}(x - 0.5) - x), \tag{52b}$$

where the Gaussians' parameters are $\sigma = 0.075$, $x_1 = 0.15$ and $x_2 = 0.85$. This initial condition is plotted in Fig. 7. All the simulations are made with a mesh of 100 cells and periodic boundary conditions are used. The coefficient $4/3$ is chosen such that m_1 is symmetric in the domain, and since u is antisymmetric, such a configuration triggers a stationary δ -shock. Remark that the monokinetic assumption (2) is not valid at

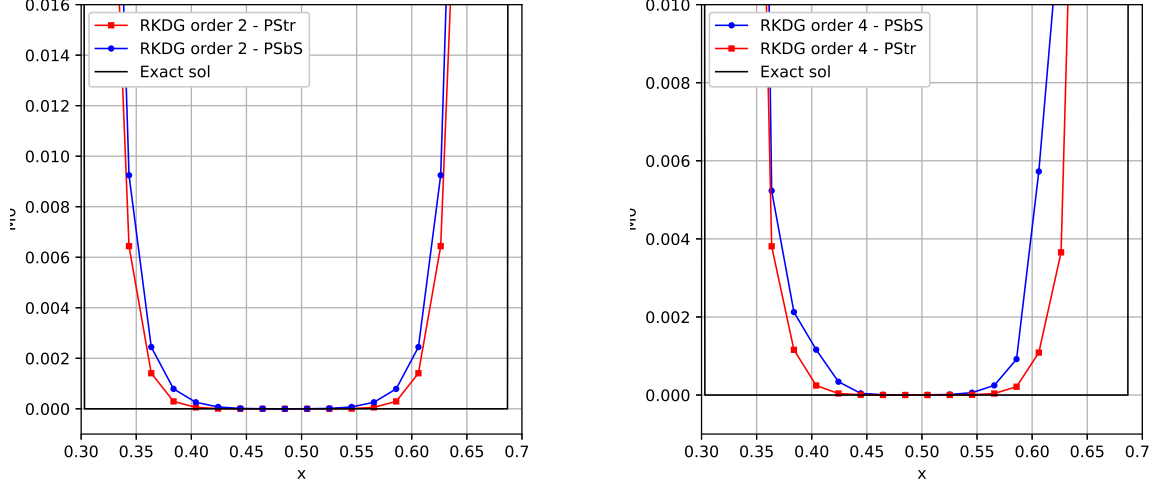


Figure 6: Zoom of the m_0 -plot in the vacuum region $x \in [0.3, 0.7]$.

this location, and the kinetic solution to (1) is composed of the sum of the two distributions defined in (52a) crossing each others at the velocity given (52b). On the contrary for the moment solution to (6), the two masses do not cross each others, but enter into a stationary δ -shock located in $x = 0.5$. This δ -shock is the sum of the masses coming from both sides of the shock which are symmetric for m_1 but not for m_0 , $m_{1/2}$ and $m_{3/2}$. Fig. 8 shows m_0 and m_1 with the different schemes at $t^N = 0.4$. The limitations mainly activate in the region where the two masses enter in the δ -shock. In this region, the moments \mathbf{m} are far from the boundary $\partial\mathcal{R}$. Therefore, the two limitations gives again identical values and only one (those with \mathcal{P}_{Str}) is shown in Fig. 8. Due to the shape of the initial velocity (see Fig. 7, right), it is crucial to impose velocity bounds that are computed locally (as in (8)) rather than globally (as in (7)) in order to filter spurious oscillations. One observes furthermore in Fig. 8 (right) that the velocity profile with the RKDG schemes is less diffused compared to the one with the KfV scheme.

Finally, we compare in Table 2 the value of the moment vector inside the δ -shocks with the exact solution. Straightforward computations leads to an exact value

$$m_\alpha^\delta = (\alpha + 1)^{-1} \left(\int_{1/6}^{1/2} G(x, x_1) dx + \frac{4}{3} (1 - 0.5^{\alpha+1}) \int_{1/2}^{5/6} G(x, x_2) dx \right).$$

The exact solution is compared to the approximate solutions given by the second and third order RKDG schemes and the second order KfV scheme. The value of the moments $m_0, m_{1/2}, m_1, m_{3/2}$ and $m_1 u$ inside the δ -shocks are presented in Table 2. The limitation strategy in the RKDG schemes uses the straight projection (45) and the KfV scheme, a minmod limiter. All the schemes overestimate the value of the moment vector inside the δ -shocks. The third order RKDG scheme gives the most accurate approximation of the exact solution compared to the two other schemes for the moments $m_0, m_{1/2}, m_1$ and $m_{3/2}$. Comparing the two second order schemes, the KfV scheme outperforms the RKDG one. As a summary, the present high-order schemes remain robust when tackling one of the singularities of the moments model based on monokinetic assumption, the δ -shocks.

4.4 2D δ -shock test case

This case corresponds to (6) with $U = (\mathbf{m}^T, q^T)^T$ where $q = (q_1, q_2)^T$ and $x = (x_1, x_2)^T$ are two-dimensional. The initial condition are:

$$m_\alpha(x, 0) = (\alpha + 1)^{-1}, \quad q(x, 0) = -0.25 \times m_1(x, 0) \times (\text{sign}(x_1), \text{sign}(x_2))^T. \quad (53)$$

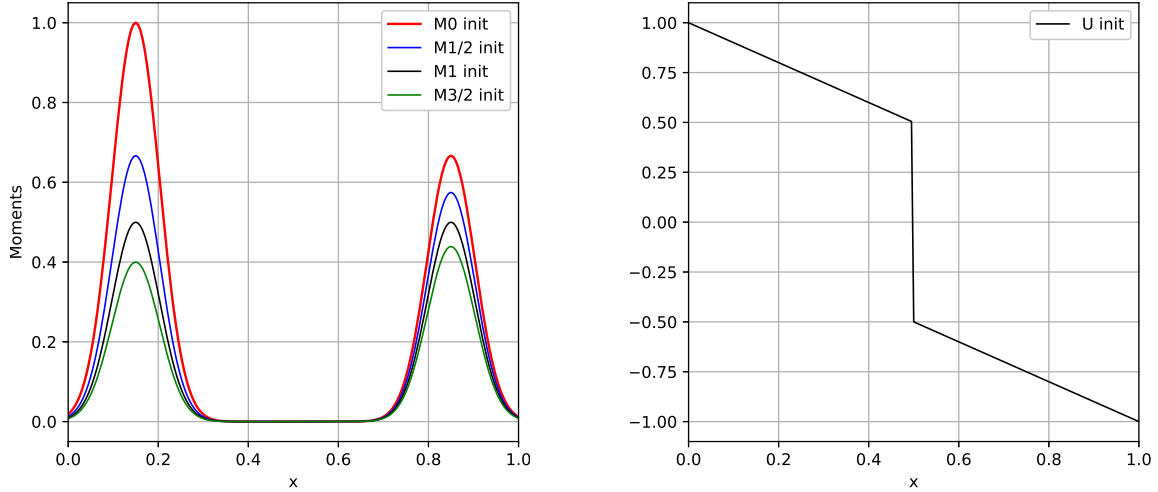


Figure 7: Initial conditions (52) on m_α (left) for $\alpha = 0, 1/2, 1, 3/2$ and on $u \equiv q/m_1$ (right).

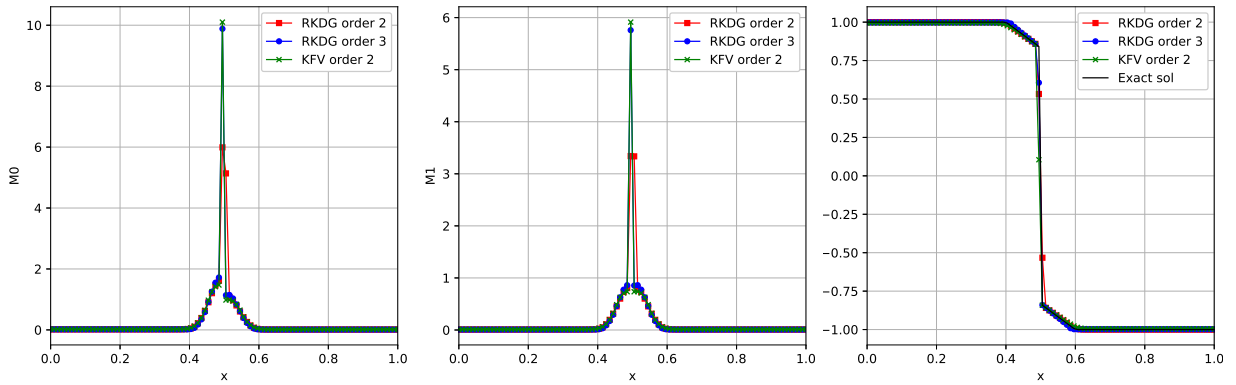


Figure 8: Moments m_0 (left) and m_1 (middle) and velocity profile $u = q/m_1$ obtained with the RKDG schemes and the KFV scheme at $t = 0.4$ with the initial conditions (52).

| | m_0 | $m_{1/2}$ | m_1 | $m_{3/2}$ | q |
|--------------|--------|-----------|--------|-----------|--------|
| Exact | 0.083 | 0.062 | 0.05 | 0.042 | 0 |
| RKDG Order 2 | 0.1112 | 0.0828 | 0.0667 | 0.056 | 0.0105 |
| RKDG Order 3 | 0.0988 | 0.0722 | 0.0575 | 0.0479 | 0.0071 |
| KFV Order 2 | 0.101 | 0.0740 | 0.0591 | 0.0493 | 0.0062 |

Table 2: Moment vector inside the δ -shock with the different schemes.

The computational domain $[-1, 1]^2$ is meshed with 100^2 uniform cells with $\Delta x = \Delta y = 1/50$. Homogeneous Dirichlet boundary conditions are imposed. The final time is $t^N = 0.5$. The initial configuration has a uniform mass m_1 over the domain, and the initial velocity profile (53) is such that δ -shocks form along the axis $x = 0$ and $y = 0$. This symmetric setup also leads to a mass concentration at the origin. Fig. 9 shows the numerical solution m_1 obtained with second order RKDG scheme with the straight projection limitation (45). The result is similar to the 1D case in Subsection 4.3. As expected, the mass accumulates along the axes and at the origin. The robustness and stability of the RKDG scheme with the limitation procedure are preserved in 2D even in presence of δ -shock singularities.

This case is modified in Section A.1 and A.2 of Appendix A with asymmetric δ -shocks or δ -shocks not

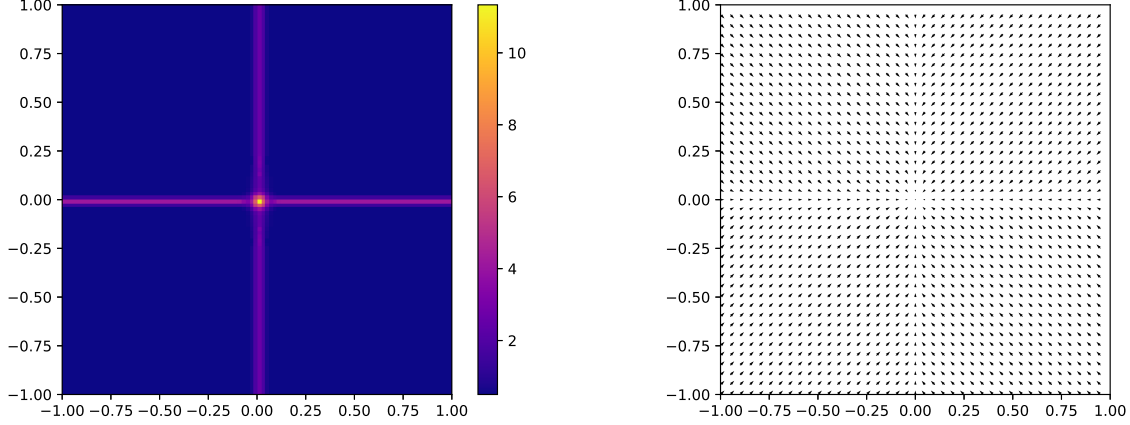


Figure 9: Moment m_1 (left) and velocity field (right) at $t^N = 0.5$ with the initial condition (53).

aligned with the mesh.

4.5 2D δ -vacuum test case

Eventually, we consider the initial condition:

$$m_\alpha(x, 0) = (\alpha + 1)^{-1}, \quad q(x, 0) = 0.4 \times m_1(x, 0) \times (\text{sign}(x_1), \text{sign}(x_2))^T. \quad (54)$$

The numerical parameters are identical to the previous case. The initial mass m_1 is uniform. The initial velocity is directed outwards the axes $x = 0$, $y = 0$ such that it generates vacuum region $m_\alpha \approx 0$ around these two axes. Fig. 10 m_0 with the second order RKDG scheme with the straight limitation (45). The

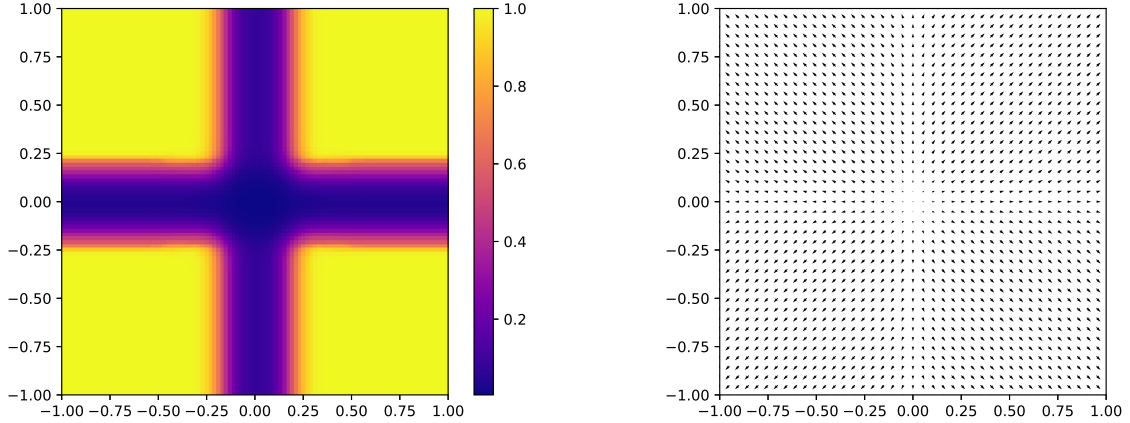


Figure 10: Moment m_0 (left) and velocity field (right) at $t^N = 0.5$ with the initial condition (54).

numerical results agree qualitatively with the expected solution. This case is modified to generate a vacuum not aligned with mesh in Section A.3 of Appendix A.

5 Conclusion

The purpose of this contribution has been to construct high-order RKDG schemes able to cope efficiently with the peculiarities of a weakly hyperbolic system of moments equations, which is frequently encountered

in fluid mechanics and combustion applications, where a spray of polydisperse droplets is to be found either coupled to a gaseous flow field or as a small scale modeling in two-scale gas-liquid flows. The key feature of the proposed method has been to maintain high-order for smooth solutions but also to cope efficiently with convex admissibility set preservation, be it the moment space or the maximum principle on velocity, in the presence of singularities and void. A specific limitation procedure has been introduced in order to limit numerical diffusion when vacuum is created and to handle properly singularities. The proposed strategy has been shown to be competitive in terms of accuracy and of robustness compared to a reference kinetic finite volume scheme [8, 34, 24] at second order but also has the ability to reach third and fourth order within the same paradigm. Such a feature is essential for applications where spray dynamics experience preferential concentration and where vacuum is always present, thus influencing the local mixture fraction and the evaporation rate for combustion applications in particular but it has a much broader range of applications. The present piece of work is also the building block for the construction of a numerical strategy for more complex multi-variate cases such as in [40] for oscillating droplets flows. The aim is to simulate a cloud of non spherical oscillating droplets by taking into account the geometrical dynamics described by the phase variables. This is work in progress.

A Additional numerical results

Here we include some additional 2D numerical results in order to illustrate the behavior of the second order RKDG scheme associated to the straight projection (45). In all the numerical experiments below, the computational domain $[-1, 1]^2$ is divided into 100×100 uniform cells with $\Delta x = \Delta y = 1/50$. Homogeneous Dirichlet boundary conditions are imposed and the computations are performed until the final time t^N . In all the cases, we solve a Riemann problem where the initial data in each quadrant are constants. The quadrants are indexed conventionally: from the first to the fourth in a counterclockwise direction, starting from the right upper one.

A.1 Asymmetric δ -shock case

The initial condition is:

$$\begin{aligned} m_\alpha(x, 0) &= (\alpha + 1)^{-1} \left(\mathbf{1}_{(\mathbb{R}^+)^2} + \mathbf{1}_{(\mathbb{R}^-)^2} + \frac{4}{3}(1 - 0.5^{\alpha+1})(\mathbf{1}_{\mathbb{R}^+ \times \mathbb{R}^-} + \mathbf{1}_{\mathbb{R}^- \times \mathbb{R}^+}) \right) (x), \\ q(x, 0) &= -0.1 \times m_1(x, 0) \times (\text{sign}(x_1), \text{sign}(x_2)). \end{aligned} \quad (55)$$

We start from the configuration where the mass m_1 is equally distributed over the domain, but not the other moments. The initial velocity profile (55) is chosen to generate stationary δ -shocks along the axes $x = 0$ and $y = 0$. Fig. 11 shows the moments m_0 and m_1 at final time $t^N = 0.25$. The moments m_0 , $m_{1/2}$ and $m_{3/2}$ are different in the quadrants next to each others. The velocity profile at the end of the simulation is the same as the one observed in the 2D δ -shock test case with equally distributed mass (53).

A.2 Diagonal δ -shock case

We use the initial condition:

$$\begin{aligned} m_\alpha(Rx, 0) &= (\alpha + 1)^{-1}, \\ q(Rx, 0) &= -0.1 \times m_1(x, 0) \times (\text{sign}(x_1), \text{sign}(x_2)). \end{aligned} \quad (56)$$

with $Rx = (\cos(\pi/3)x_1 + \sin(\pi/3)x_2, -\sin(\pi/3)x_1 + \cos(\pi/3)x_2)$. The initial mass m_1 is equally distributed over the domain and the initial velocity (56) is chosen to generate δ -shocks along the diagonals $y = x\sqrt{3}$ and $y = x/\sqrt{3}$. The velocity is heading towards the two diagonals and the solution of the moment system (6) yields δ -singularities located at the origin and the two axes. This is illustrated on Fig. 12 which shows m_1 and $u \equiv q/m$ at final time $t^N = 0.25$. The numerical solution exhibits the expected behavior, i.e. stationary δ -shocks along the two diagonals $y = x\sqrt{3}$ and $y = x/\sqrt{3}$ (see Fig. 12 (left) for moment m_1). The result

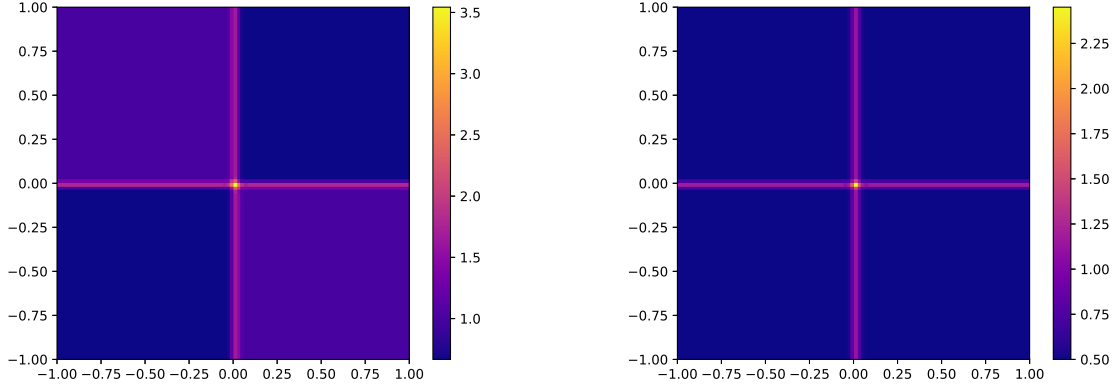


Figure 11: Moments m_0 (left) and m_1 (right) at $t^N = 0.25$ with the initial condition (55).

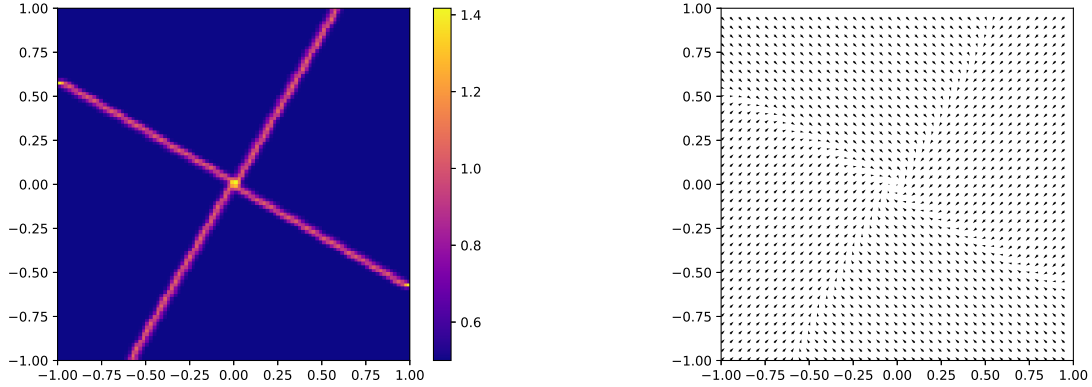


Figure 12: Moment m_1 (left) and velocity (right) at $t^N = 0.25$ with the initial condition (56).

is similar to the 2D case (53) after applying a rotation of angle $\pi/3$. The robustness and stability of the RKDG scheme and its limitation procedure are preserved even in the case of δ -shock singularities that are not aligned with the mesh.

A.3 Diagonal vacuum case

Eventually, we consider the initial condition:

$$m_\alpha(Rx, 0) = (\alpha + 1)^{-1}, \quad q(Rx, 0) = 0.4 \times m_1(x, 0) \times (\text{sign}(x_1), \text{sign}(x_2)). \quad (57)$$

The initial mass m_1 is equally distributed over the domain and the initial velocity (57) generates vacuum $m_\alpha \approx 0$ along the diagonals $y = x\sqrt{3}$ and $y = x/\sqrt{3}$. Indeed, the velocity is heading outwards these axes. Fig. 13 shows the moment m_1 and the velocity $u \equiv q/m_1$ at final time $t^N = 0.5$. Again, the behavior of the numerical method agree qualitatively with the structure of the expected solution. Indeed, we can observe that the numerical solution approximates well the vacuum that are not aligned with the mesh. Because of the presence of vacuum, the limitation procedure is activated in order to avoid non-realizable vector of moments as in the 2D case of (54) after applying a rotation of angle $\pi/3$.

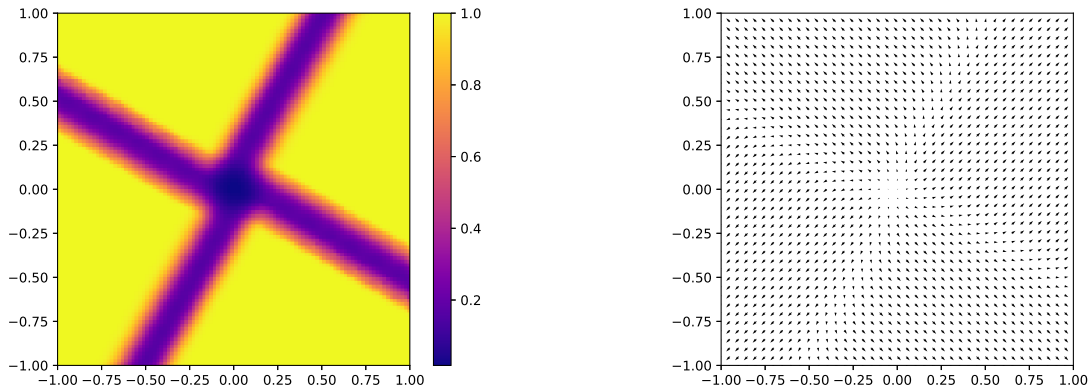


Figure 13: Moment m_0 (left) and velocity (right) at $t^N = 0.5$ with the initial condition (57).

References

- [1] N. I. Akhiezer. *The classical moment problem*. Edinburgh : Oliver & Boyd, 1965.
- [2] R. Anderson, V. Dobrev, T. Kolev, D. Kuzmin, M. Quezada de Luna, R. Rieben, and V. Tomov. High-order local maximum principle preserving MPP discontinuous Galerkin finite element method for the transport equation. *J. Comput. Phys.*, 334:102–124, 2017.
- [3] F. Bouchut. On zero pressure gas dynamics. In *Advances in kinetic theory and computing*, volume 22 of *Ser. Adv. Math. Appl. Sci.*, pages 171–190. World Sci. Publ., 1994.
- [4] F. Bouchut and F. James. One-dimensional transport equations with discontinuous coefficients. *Non-linear Anal. TMA*, 32(7):891–933, 1998.
- [5] F. Bouchut and F. James. Solutions en dualité pour les gaz sans pression. *C. R. Acad. Sci. Paris, Série 1*, pages 1073–1078, 1998.
- [6] F. Bouchut and F. James. Duality solutions for pressureless gases, monotone scalar conservation laws, and uniqueness. *Commun. Partial Diff. Eq.*, 24:2173–2189, 1999.
- [7] F. Bouchut, F. James, and S. Mancini. Uniqueness and weak stability for multi-dimensional transport equations with one-sided Lipschitz coefficient. *Ann. Scuola Normale Superiore di Pisa*, 4(1):1–25, 2005.
- [8] F. Bouchut, S. Jin, and X. Li. Numerical approximations of pressureless and isothermal gas dynamics. *SIAM J. Num. Anal.*, 41:135–158, 2003.
- [9] Y. Brenier and E. Grenier. Sticky particles and scalar conservation laws. *SIAM J. Numer. Anal.*, 35(6):2317–2328, 1998.
- [10] Z. Cai. On the finite volume element method. *Numer. Math.*, 58:713–735, 1990.
- [11] Z. Cai. *Handbook of Numerical Analysis*, volume 7. Elsevier, 2000.
- [12] C. Chalons, D. Kah, and M. Massot. Beyond pressureless gas dynamics: quadrature-based velocity moment models. *Commun. Math. Sci.*, 10(4):1241–1272, 2012.
- [13] Nattaporn Chuenjarern, Ziyao Xu, and Yang Yang. High-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes. *J. Comput. Phys.*, 378:110–128, 2019.

- [14] B. Cockburn and C.-W. Shu. The Runge-Kutta Discontinuous Galerkin Method for Conservation Laws V. *J. Comput. Phys.*, 141(2):199–224, 1998.
- [15] B. Cockburn and C.-W. Shu. *Discontinuous Galerkin Methods: Theory, computation and applications*. Springer, 2000.
- [16] Bernardo Cockburn and Chi-Wang Shu. TVB Runge-Kutta Local Projection Discontinuous Galerkin Finite Element Method for Conservation Laws II: General Framework. *Math. Comp.*, 52(186):411–435, 1989.
- [17] P. Cordesse. *Contribution to the study of combustion instabilities in cryotechnic rocket engines : coupling diffuse interface models with kinetic-based moment methods for primary atomization simulations*. PhD thesis, CentraleSupélec, 2020.
- [18] R. Curto and L. A. Fialkow. Recursiveness, positivity, and truncated moment problems. *Houston J. Math.*, 17(4):603–634, 1991.
- [19] H. Dette and W. J. Studden. *The theory of canonical moments with applications in statistics, probability, and analysis*. John Wiley & Sons Inc., 1997.
- [20] D. Di Pietro and A. Ern. *Mathematical Aspects of Discontinuous Galerkin Methods*, volume 69 of *SMAI Mathématiques et Applications*. Springer, 2012.
- [21] D. A. Drew. Evolution of geometric statistics. *SIAM J. Appl. Math.*, 50(3):649–666, 1990.
- [22] F. Druil. *Modélisation et simulation Eulériennes des écoulements diphasiques à phases séparées et dispersées : développement d’une modélisation unifiée et de méthodes numériques adaptées au calcul massivement parallèle*. PhD thesis, Université Paris-Saclay, 2017.
- [23] M. Essadki. *Contribution to a unified Eulerian modeling of fuel injection: from dense liquid to polydisperse evaporating spray*. PhD thesis, CentraleSupélec, 2018.
- [24] Mohamed Essadki, Stéphane de Chaisemartin, Frédérique Laurent, and Marc Massot. High-order moment model for polydisperse evaporating sprays towards interfacial geometry. *SIAM J. Appl. Math.*, 78(4):2003–2027, 2018.
- [25] Mohamed Essadki, Stéphane de Chaisemartin, Marc Massot, Frédérique Laurent, Adam Larat, and Stéphane Jay. Adaptive Mesh Refinement and High-Order Geometrical Moment Method for the Simulation of Polydisperse Evaporating Sprays. *Oil & Gas Sci. Tech. - Rev. IFP Energies nouvelles*, 71(5), 2016.
- [26] A. F. Filippov. *Differential Equations with Discontinuous Righthand Sides*. Springer, 1988.
- [27] Sigal Gottlieb and Chi-Wang Shu. Total variation diminishing runge-kutta schemes. *Math. Comp.*, 67:73–85, 08 1996.
- [28] Sigal Gottlieb, Chi-Wang Shu, and Eitan Tadmor. Strong stability-preserving high-order time discretization methods. *SIAM Rev.*, 43(1):89–112, 2001.
- [29] J.-L. Guermond, B. Popov, and I. Tomas. Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems. *Comput. Meth. Appl. Mech. Eng.*, 347:143–175, 2019.
- [30] H. Hajduk. Monolithic convex limiting in discontinuous Galerkin discretizations of hyperbolic conservation laws. *Computers Math. Appl.*, 87:120–138, 2021.
- [31] A. Harten. High resolution schemes for hyperbolic conservation laws. *J. Comput. Phys.*, 49(3):357–393, 1983.

- [32] F. Hausdorff. Summationmethoden und momentfolgen. *Math. Z.*, (9):74–109, 1921.
- [33] Jan S. Hesthaven. *Numerical Methods for Conservation Laws*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2018.
- [34] D. Kah, F. Laurent, M. Massot, and S. Jay. A high-order moment method simulating evaporation and advection of a polydisperse spray. *J. Comput. Phys.*, 231(2):394–422, 2012.
- [35] H. C. Kranzer and B. L. Keyfitz. A strictly hyperbolic system of conservation laws admitting singular shocks. In *Nonlinear Evolution Equations that Change Type*, volume 27, pages 107–125. IMA Volumes in Mathematics and its Applications, 1990.
- [36] J.-B. Lasserre. *Moment, positive polynomials, and their applications*, volume 1. Imperial college press, 2009.
- [37] Ph. Le Floch. An existence and uniqueness result for two nonstrictly hyperbolic systems. In *Nonlinear Evolution Equations that Change Type*, volume 27, pages 126–138. IMA Volumes in Mathematics and its Applications, 1990.
- [38] Yimin Lin, Jesse Chan, and Ignacio Tomas. A positivity preserving strategy for entropy stable discontinuous Galerkin discretizations of the compressible Euler and Navier-Stokes equations. *J. Comput. Phys.*, 475:111850, 2023.
- [39] A. Loison. *Unified two-scale Eulerian multi-fluid modeling of separated and dispersed two-phase flows*. PhD thesis, IP Paris, 2024.
- [40] A. Loison, S. Kokh, M. Massot, and T. Pichard. Two-phase flow reduced-order model: a two-scale model of polydisperse oscillating droplets. *Submitted, arXiv:2308.15641*, 2023.
- [41] Yu Lv, Yee Chee See, and Matthias Ihme. An entropy-residual shock detector for solving conservation laws using high-order discontinuous Galerkin methods. *J. Comput. Phys.*, 322:448–472, 2016.
- [42] D. Marchisio and R. Fox. Solution of population balance equations using the direct quadrature method of moments. *J. Aerosol Sci.*, 36:43–73, 2005.
- [43] M. Massot, F. Laurent, S. de Chaisemartin, L. Fréret, and D. Kah. Eulerian multi-fluid models: modeling and numerical methods. In *Modelling and Computation of Nanoparticles in Fluid Flows*, Lectures Notes of the von Karman Institute. NATO RTO-EN-AVT-169, 2009. Available at <http://hal.archives-ouvertes.fr/hal-00423031/en/>.
- [44] M. Massot, F. Laurent, D. Kah, and S. de Chaisemartin. A robust moment method for evaluation of the disappearance rate of evaporating sprays. *SIAM J. Appl. Math.*, 70:3203–3234, 2010.
- [45] A. Mazaheri, C.-W. Shu, and V. Perrier. Bounded and compact weighted essentially nonoscillatory limiters for discontinuous Galerkin schemes: Triangular elements. *J. Comput. Phys.*, 395:461–488, 2019.
- [46] R. McGraw. Description of aerosol dynamics by the quadrature method of moments. *Aerosol Sci. Tech.*, (27):255–265, 1987.
- [47] W. Pazner. Sparse invariant domain preserving discontinuous Galerkin methods with subcell convex limiting. *Comput. Meth. Appl. Mech. Eng.*, 382:113876, 2021.
- [48] T. Pichard. A moment closure based on a projection on boundary of the realizability domain: 1d case. *Kin. rel. mod.*, 13(6):1243–1280, 2020.
- [49] S. B. Pope. The evolution of surfaces in turbulence. *Int. J. Engng Sci.*, 26(5):445–269, 1988.

- [50] F. Poupaud and M. Rascole. Measure solutions to the linear multi-dimensional transport equation with non-smooth coefficients. *Commun. Partial Diff. Eq.*, 22(1-2):225–267, 1997.
- [51] Jianxian Qiu and Chi-Wang Shu. A comparison of troubled-cell indicators for runge–kutta discontinuous galerkin methods using weighted essentially nonoscillatory limiters. *SIAM J. Sci. Comput.*, 27(3):995–1013, 2005.
- [52] Jianxian Qiu and Chi-Wang Shu. Runge-Kutta discontinuous Galerkin method using WENO limiters. *SIAM J. Sci. Comput.*, 26(3):907 – 929, 2005.
- [53] Andrés M Rueda-Ramírez, Benjamin Bolm, Dmitri Kuzmin, and Gregor J Gassner. Monolithic Convex Limiting for Legendre-Gauss-Lobatto Discontinuous Galerkin Spectral Element Methods, 2023.
- [54] K. Schmüdgen. *The moment problem*. Springer, 2018.
- [55] F. Schneider. Kershaw closures for linear transport equations in slab geometry ii: High-order realizability-preserving discontinuous-galerkin schemes. *J. Comput. Phys.*, 322, 02 2016.
- [56] C.-W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock-capturing schemes. *J. Comput. Phys.*, 77:439–471, 1988.
- [57] D. Tan, T. Zhang, and Y. Zheng. Delta-shock waves as limits of vanishing viscosity for hyperbolic systems of conservation laws. *J. Diff. Eq.*, 112:1–32, 1994.
- [58] F. A. Williams. Spray combustion and atomization. *Phys. Fluids*, 1:541–545, 1958.
- [59] Y. Xing, X. Zhang, and C.-W. Shu. Positivity-preserving high-order well-balanced discontinuous Galerkin methods for the shallow water equations. *Adv. Water Resources*, 33, 12 2010.
- [60] Y. Yang, D. Wei, and C.-W. Shu. Discontinuous Galerkin method for Krause’s consensus models and pressureless Euler equations. *J. Comput. Phys.*, pages 109–127, 2013.
- [61] X. Zhang. *Maximum-Principle-Satisfying and Positivity-Preserving High-Order Schemes for Conservation Laws*. PhD thesis, Brown, 2011.
- [62] X. Zhang. On positivity-preserving high-order discontinuous Galerkin schemes for compressible Navier–Stokes equations. *J. Comput. Phys.*, 328:301–343, 2017.
- [63] X. Zhang and C.-W. Shu. On maximum-principle-satisfying high-order schemes for scalar conservation laws. *J. Comput. Phys.*, 229(9):3091–3120, 2010.
- [64] X. Zhang and C.-W. Shu. On positivity-preserving high-order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes. *J. Comput. Phys.*, 229:8918–8934, 2010.
- [65] X. Zhang and C.-W. Shu. Positivity-preserving high-order discontinuous Galerkin schemes for compressible Euler equations with source. *J. Comput. Phys.*, 230(4):1238–1248, 2011.
- [66] Xiangxiong Zhang and Chi-Wang Shu. Maximum-principle-satisfying and positivity-preserving high-order schemes for conservation laws: Survey and new developments. *Proc. Royal Soc. A*, 467, 05 2011.
- [67] Xinghui Zhong and Chi-Wang Shu. A simple weighted essentially nonoscillatory limiter for runge-kutta discontinuous galerkin methods. *J. Comput. Phys.*, 232(1):397–415, 2013.