

POSITIVITY-PRESERVING DISCONTINUOUS GALERKIN SCHEMES FOR LINEAR VLASOV-BOLTZMANN TRANSPORT EQUATIONS

YINGDA CHENG ^{*}, IRENE M. GAMBA [†], AND JENNIFER PROFT [‡]

Abstract. We develop a high-order positivity-preserving discontinuous Galerkin (DG) scheme for linear Vlasov-Boltzmann transport equations (Vlasov-BTE) under the action of quadratically confined electrostatic potentials. The solutions of such BTEs are positive probability distribution functions and it is very challenging to have a mass-conservative, high-order accurate scheme that preserves positivity of the numerical solutions in high dimensions. Our work extends the maximum-principle-satisfying scheme for scalar conservation laws by X. Zhang and C.-W. Shu [52] to include the linear Boltzmann collision term. The DG schemes we developed conserve mass and preserve the positivity of the solution without sacrificing accuracy. A discussion of the standard semi-discrete DG schemes for the BTE are included as a foundation for the stability and error estimates for this new scheme. Numerical results of the relaxation models are provided to validate the method.

Key words. Boltzmann transport equations, discontinuous Galerkin finite element methods, positivity-preserving schemes, stability, error estimates, relaxation models.

AMS subject classifications. 65M60, 76P05, 74S05

1. Introduction. The primary behavior of kinetic models is driven by the interaction between transport and collisional operators. They describe the evolution of a probability density mass associated to interacting particle systems, where the transport is modeled by a first order linear differential operator along a particle path modified by the presence of an electric field which is balanced by a particle interacting integral (collisional) operator, of dissipative nature, modeling the rate of gain and loss of the probability density mass due to the interactions. Such integro-differential operator is referred to as a Vlasov-Boltzmann Transport equation (Vlasov-BTE).

In this particular manuscript we focus on linear collisional structures where the transport along parabolic energy bands in phase space are modified by an external electric field. Because of the mathematical difficulties arising by this interaction operators as well as the large number of unknowns, the numerical approximation of such kinetic transport models is highly demanding from a computational standpoint. Traditional computational schemes for this type of linear transport developed several decades ago were based in Discrete Simulations Monte Carlo (DSMC) [3, 39] or more recently deterministic approaches initiated by [26, 37] and more recently extended WENO finite difference based schemes [6, 7, 8, 9, 5]

The discontinuous Galerkin (DG) finite element method, which is well suited for *hp*-adaptivity and parallel implementation, recently gained growing attentions in this field [12, 10, 13].

In this particular manuscript we focus on numerical analytical issues as well as positivity preserving schemes by DG approach for these type transport problems with degenerate linear collisional forms. These methods were originally developed for computing hyperbolic conservation laws. These DG methods are compact, locally conservative and high-order accurate. They can be easily extended to very high

^{*}Department of Mathematics and ICES, University of Texas at Austin Austin, TX 78712 U.S.A.
ycheng@math.utexas.edu

[†]Department of Mathematics and ICES, University of Texas at Austin Austin, TX 78712 U.S.A.
gamba@math.utexas.edu

[‡]ICES, University of Texas at Austin Austin, TX 78712 U.S.A. jennifer@ices.utexas.edu

order and be adapted to unstructured grids. More specifically, we prove analytical properties of the approximating scheme such as numerical stability and error estimates and develop a positivity-preserving DG scheme that ensures the positivity of numerical solution and is genuinely high order for arbitrary order of discretization.

The Vlasov-BTE describes the evolution properties of a probability distribution function (*pdf*) $f = f(t, x, v)$ representing the probability of finding a particle at time t with position at x and phase velocity v . This model represents a dilute or rarefied gaseous state corresponding to a probabilistic description of particles when the transport is given by a classical Hamiltonian with acceleration component given by the action of a Lorentzian force and particle interactions taken into account in a dissipative, non-local framework modeled by the so called collision operator.

In the case of charged transport in the absence of strong magnetic effects, the Lorentzian force field is reduced to an electric force field corresponding to an electrostatic potential modeled by the Poisson equation for total charges, as in the mean field theory approximation. The collisional integral, classically modeled by bilinear non-local interactions, satisfy the Pauli Exclusion Principle for which Fermi-Dirac distributions are in their null spaces. However, it is well known [44, 48] that for low density or hot temperature transport the collisional form is well approximated by a linear non-local operator modeling charged electrons interacting with background impurities or other carrier type.

In our framework, such linear collision operator, denoted by $Q(f, \sigma)$, models a gain/loss of probability rates, where the function σ represents transition probabilities of scatters from one state to another (scattering mechanisms) multiplied by given distributions.

We consider then an initial value problem to the Vlasov-BTE equation, which has a linear hyperbolic component of transport and a dissipative mechanism due to the collision operator. Its solutions are positive probability distributions, so-called Borel measures. In the absence of boundary injection, they conserve the total mass, $\int_{\mathbb{R}^d_x} \int_{\mathbb{R}^d_v} f(t, x, v) dx dv$, at any time.

The development and analytical properties of any discrete scheme to compute the Vlasov-BTE deterministically requires that the approximation is done in a finite domain, both in physical x -space and phase v -space, even though the setting of the problem is in all \mathbb{R}^d in v -space. In particular, the choice of the finite computational domain depends very much on the nature of the problem to be solved and analyzed. More specifically, since the computed problem admits a stationary state, one may use an educated guess for the choice of the computational domain for which the stationary state has the same mass as that of the initial state.

However, if the solution of the problem does not stabilize, due for instance, to external unstable forcing or a bifurcation phenomena, one must carefully monitor the evolution of the *pdf* and secure that at any time step there is no significant loss of mass with respect to the initial state and its discretization in the corresponding truncated domain. This procedure could be done by an *a posteriori* algorithm using data postprocessing but is beyond the scope of this paper.

Only in recent years has the Boltzmann equation been tackled numerically with particular attention paid to accuracy and computational costs. The mathematical difficulties related to the Boltzmann equation make it extremely difficult, if not impossible, to determine analytic solutions in most physically relevant situations. Consequently, numerical methods, particularly highly accurate deterministic methods, must be used to obtain valid approximations of the solution.

We mentioned that the well known Discrete Simulation Monte-Carlo (DSMC) methods have traditionally been used to numerically model classical kinetic equations in rarefied gas dynamics [3, 39], and in charge transport in submicron structures [48, 1]. Although Monte-Carlo techniques guarantee efficiency and preservation of the some fundamental physical properties (ex. observables), it is well established that in the presence of non-equilibrium stationary states or near continuum regimes, avoiding statistical fluctuations in the resulting solution may become extremely expensive. Consequently, deterministic methods are more competitive when greater accuracy is required. A comparative study of DSMC and deterministic solvers by WENO [8, 9] and DG schemes [11, 10, 12, 13, 14] for short base channels models by 1-dimensional electron transport models and 2-dimensional MESFET and MOSFET devices have been carried out. They are important discrepancies between both methods near ohmic contacts and interface layers. It is also well accepted that that a more refined Monte Carlo code is able to generate better hydrodynamical profiles the authors observed an excellent agreement for grid points where the solution shows strong gradients, better resolution and implementations of boundary conditions as well as comparative computational time by mesh adaptivity. Other popular deterministic solvers for Boltzmann equations have been developed in the last two decades. They include the discrete velocity methods initiated by [4], spectral methods [40, 27], and finite difference schemes [26, 8, 9] among many others. However to the best of our knowledge neither discrete velocity models nor spectral methods, which have been essentially computational tools for non-linear collisional operators, have been implemented in the modeling of charged electron transport in nano devices where the collisional forms are taken to be linear in the hot electron or low density transport.

It is quite interesting to observe that the first DG method approach, introduced first in 1973 by Reed and Hill [43] for neutron transport, are truly relevant to our studies. Later, Lesaint and Raviart [36] performed the first convergence study for the original DG method. Cockburn *et al.* in a series of papers [22, 21, 20, 18, 23] developed the Runge-Kutta DG (RKDG) method for hyperbolic equations. For more details about the RKDG scheme, one can refer to the survey paper [24] and the references within. Very recently in [11, 10, 12, 13], DG schemes for solving the Boltzmann-Poisson system in semiconductor device modeling are proposed. The scheme shows excellent behavior in terms of quality of the numerical solution and computational efficiency. The DG schemes provide a tool to accurate transient calculations of strong non-equilibrium stationary states, and allow for the explicit computation of the probability density function and, consequently, all moments with greater accuracy as compared to Monte Carlo methods. Moreover, their inherent flexibility allows for the resolution of temporal and spatial non-uniformities. The additional computational costs associated with deterministic methods can be greatly reduced by focusing on regions of interest, and coarsening the mesh in areas such as the “tails” of the probability density function where there is little action in the behavior of the solution.

In this manuscript, we propose a positivity preserving DG scheme and rigourously analyze the computational scheme in the case where the solution corresponds to the initial value problem of the linear BTE evolving under the field generated by a confined potential. In this case it is easy to find stationary states for relaxation models and their corresponding decay rates to equilibrium [33, 38]. We will use these properties to select an appropriate truncation for the computational domain. L^2 stability and error estimates for general order of approximations are provided. Since the solutions to the BTEs are positive *pdfs*, it is desired for the numerical scheme to generate

positive solutions as approximated negative *pdf* not only are non-physical, but also will produce negative contribution to higher order moments such as energy or heat flow. So it is highly desired to have mechanisms for high order schemes that preserve the positivity of the solution.

We recall such numerical positivity property is easily verified for the piecewise constant DG schemes. However, typical DG or finite difference schemes of order higher than one are no longer monotone and will not preserve the positivity of the solution. In addition monotone schemes can only be first order accurate and total variation diminishing (TVD) and those that satisfy the maximum principle are at most second order accurate in the L^1 norm.

It is, therefore, a difficult task to preserve the positivity and high-order accuracy of the solution at the same time. This is a drawback for most high order numerical methods.

Our current work uses a maximum-principle-satisfying limiter that has been recently proposed by Zhang *et al.* in [52] for conservation laws and develop it for the Boltzmann equation. This method has been used to develop positivity-preserving schemes for compressible Euler [53] and shallow water equations [55]. The resulting scheme preserves positivity of numerical solution in cell average sense and is genuinely high order accurate for arbitrary order of discretization. In addition, the limiter is applied as a post-processing step, which induces minimal computational cost. In particular, most of the new work presented here focus on the treatment of the approximation of the collisional integral for the preservation of positivity under the proposed new scheme.

It is to our conclusion that, based on the analysis of the standard DG methods and the positivity-preserving property of this modified scheme, we have obtained a highly-efficient deterministic solver that obtains accurate and physically relevant solutions to the linear BTEs.

The rest of the paper is organized as follows. In Section 1, we introduce the Vlasov-BTE and recent global regularity and decay results. A discussion of the choice of computational domain is included. We summarize some important properties of the collisional operator in Section 2. The traditional semi-discrete DG scheme and the positivity-preserving DG scheme will be formulated in Section 3. Those two schemes are analyzed in Sections 4 and 5 respectively. Numerical results for relaxation models are given in Section 6. Finally, concluding remarks and remarks on future work are provided in Section 7.

We consider the initial value problem for the linear Vlasov Boltzmann Transport Equation for a charged particle distribution function $f = f(t, x, v)$ defined on $\mathbb{R}_t^+ \times \mathbb{R}_x^d \times \mathbb{R}_v^d$,

$$\begin{aligned} \frac{\partial f}{\partial t} + v f_x - \frac{e}{m} E(t, x) f_v &= Q_\sigma(f)(t, x, v), \\ f(0, x, v) &= f_0(x, v) \quad \text{in } \mathbb{R}_x^d \times \mathbb{R}_v^d \end{aligned} \quad (1.1)$$

where e is the space charge constant associated to the particle species and m is the effective particle mass constant. Denoting by $f' = f(t, x, v')$, the linear collision operator $Q_\sigma : f(t, x, v) \rightarrow Q_\sigma(f)(t, x, v)$ is given by

$$Q_\sigma(f)(t, x, v) = \int_{v' \in \mathbb{R}^d} (\sigma(x, v, v') f' - \sigma(x, v', v) f) dv', \quad (1.2)$$

where the scattering function $\sigma(x, v, v')$ is positive and may satisfy the *detailed balance principle*, usually modeled by

$$\sigma(x, v, v') = k(x, v, v') M(v), \quad \sigma(x, v', v) = k(x, v', v) M(v'). \quad (1.3)$$

Here $k(x, v, v')$ is symmetric in (v, v') and represents the transition probability of scatters passing from a state with velocity v into v' including spacial variations, and $M(v)$ is a stationary probability distribution, independent of space. Such operator models a linear degenerate collisional form, sometimes referred as a thermostat, such that the only function in its kernel is $M(v)$. When one assumes this stationary probability distribution to be an absolute Maxwellian distribution denoted by

$$\mu_\infty(v) = M(v) = \frac{\exp(-|v|^2/2\theta)}{(2\pi\theta)^{d/2}}, \quad (1.4)$$

then the collisional integral interactions with a Gaussian distributed background impurities with constant kinetic temperature θ .

In general, provided that both $M(v), |v|^2 M(v) \in L^1(\mathbb{R}^d)$, such collisional forms preserves mass, but does not conserve any other velocity moment of f . In particular it conserves neither momentum $\int_{v \in \mathbb{R}^d} v f(t, x, v) dv$ nor energy $\int_{v \in \mathbb{R}^d} |v|^2 f(t, x, v) dv$.

We also need to use that the scattering function is not singular, and so we shall use that $\int_{\mathbb{R}^d} |\sigma(x, v', v)| dv' \leq K(x, v)$ with $K(x, v)$ a function with at most linear growth (see Property 1 in the next section). Usually $\nu(x, v) = \int_{\mathbb{R}^d} |\sigma(x, v', v)| dv'$ is referred to as the collision frequency term. In particular, taking the transition probability k constant, one obtains the well establish *linear relaxation model*

$$L(f)(t, x, v) = \frac{\mu_\infty(v) \rho(t, x) - f(t, x, v)}{\tau}, \quad (1.5)$$

where $\tau = \frac{1}{k}$, and

$$\rho(t, x) = \int_{\mathbb{R}^d} f(t, x, v) dv$$

denotes the macroscopic density.

We mention that in the mean field theory approximation that accounts for long range interactions this Vlasov-Boltzmann transport equation is coupled to the Poisson equation modeling electrostatic potential due to total space charges

$$E(t, x) = -\nabla_x V, \quad \operatorname{div}_x(\epsilon_0 \nabla_x V) = e(\rho(t, x) - C(x)), \quad (1.6)$$

where ϵ_0 is the permittivity of the medium and $C(x)$ is the background distribution density. The resulting non-linear system is referred as the Boltzmann-Poisson systems that appears in the modeling of charged transport in sub-micron devices. In this case, the scattering function, derived from the Fermi's Golden Rule, may be modeled by singular distribution measure supported on energy levels associated to the underlying model of electronic band structure for the semiconducting material. A DG implementation for this systems in realistic sub-micron highly heterogeneous structures was developed by the authors in collaboration with A. Majorana and C.W. Shu, see [11, 10, 12, 13].

In the case of a smooth, bounded spatial domain $\mathcal{D}_x^d = [0, L]^d \subset \mathbb{R}_x^d$, boundary and initial conditions supplementing the Boltzmann equation (1.1) are

$$\mathcal{B}f(t, x, v) \big|_{\partial \mathcal{D}_x^d} \text{ is prescribed on } v \cdot \nu(x) < 0, \quad (1.7)$$

where $\nu(x)$ is the outer unit normal to $\partial \mathcal{D}_x^d$ and $\mathcal{B}f$ is a boundary operator. These spatial boundary conditions can be of different nature, such as particle injection, specular reflection or diffusive conditions, as well as periodic spatial conditions. In addition, our analysis will use zero particle density at the boundary essentially due to the stability nature of the problem for which the numerical analysis of the approximation is performed. This latter one is a delicate issue related to a suitable cut-off domain that we will carefully address below.

The initial condition for the particle distribution function is given by

$$f(0, x, v) = f_0(x, v), \quad \text{for } x \in \mathcal{D}_x^d, v \in \mathbb{R}_v^d, \quad (1.8)$$

in an adequate space for which existence, uniqueness and functional and decay estimates hold. We will describe the available analytical properties for the initial boundary value problem in subsection 1.1 below.

Without loss of generality, we redefine $\tilde{V} = \frac{e}{m}V$, and eliminate the tilde notation and we rewrite the transport equation (1.1) as follows

$$\frac{\partial f}{\partial t} + \alpha \cdot \nabla f = Q_\sigma(f), \quad (1.9)$$

where

$$\alpha(x, v) = \begin{pmatrix} v \\ -E(t, x) \end{pmatrix}, \quad E = -\nabla_x V, \quad \nabla = \begin{pmatrix} \nabla_x \\ \nabla_v \end{pmatrix}, \quad |\alpha| = |v| + |\nabla_x V|. \quad (1.10)$$

Clearly, $\operatorname{div}_{x,v} \alpha = 0$. The quantity $|\alpha|$ is the l_1 -distance of the vector α to the origin in \mathbb{R}^{2d} .

In order to properly truncate the computational domain, we first recall properties of existence, uniqueness and stability of the linear collisional model in a properly selected Banach space for which decay rates to equilibrium have been established. These properties allows us to conclude that, in within our truncated computational domain, the proven error estimates are going to be optimal up to a uniform in time error depending only on the total mass of the stationary state.

1.1. Notation and recent global regularity and decay results. We recall the following weighted functional spaces to be used in our estimates. Let

$$L_m^p(\Omega) \equiv \left\{ \psi : \int_{\Omega} |\psi|^p (1 + |x|^2 + |v|^2)^{m/2} dv dx < \infty \right\}$$

denote the space of L^p -integrable functions with bounded polynomial decay, both x and v , with the natural norm when $\Omega = \mathbb{R}^{2d}$ given by

$$\|\psi\|_{L_m^p(\mathbb{R}^{2d})} = \left(\int_{\mathbb{R}_x^d} \int_{\mathbb{R}_v^d} |\psi|^p (1 + |x|^2 + |v|^2)^{m/2} dv dx \right)^{1/p}. \quad (1.11)$$

Note that L_0^p is the classical L^p -space.

In the case where the force field is given externally by the gradient of a time independent potential, i.e. $E(x) = -\nabla V(x)$, where the potential $V(x)$ satisfies that its derivatives of order higher or equal than second order are bounded and that $\exp(V(x)) \in L^1$, the unique stationary state solution to the system is the global Maxwellian distribution $\mu_\infty(v)$ in v -space, defined in (1.4) (which lays in the kernel of the collision operator) multiplied by the stationary macroscopic density given by the spatial Maxwellian

$$\rho_\infty = \frac{e^{V/\theta}}{\int e^{V/\theta} dx}.$$

Consequently, the unique stationary state is given by

$$\mathcal{M}(x, v) = \rho_\infty(x) \mu_\infty(v) = \frac{e^{-(\frac{|v|^2}{2} - V(x))/\theta}}{(2\pi\theta)^{d/2} \int e^{V(x)/\theta} dx}. \quad (1.12)$$

Additionally there exist positive constants $K_1, K_2, R > 0$, such that

$$K_1|x| \leq |\nabla_x V| \leq K_2|x|, \text{ for any } |x| > R. \quad (1.13)$$

Thus, any other steady states in $\mathcal{S}'(\mathbb{R}^{2d})$ are proportional to the Maxwellian \mathcal{M} , where $\mathcal{S}'(\mathbb{R}^{2d})$ is the dual of the Schwartz class of rapidly decaying functions in \mathbb{R}^{2d} (see [33]). Without loss of generality, we take $\theta = 1$.

Finally, we mention that the collision operator $Q_\sigma(f)$ is mass preserving, positivity preserving and dissipative in the sense of non-negative in $L^1(\mathbb{R}^{2d})$. In particular it is easy to see that $v\partial_x + \partial_x V(x)\partial_v + Q$ is also mass and positivity preserving and dissipative in $L^1(\mathbb{R}^{2d})$.

In addition the initial value problem is also solvable in $L^2(\mathbb{R}^{2d})$ since the monotonicity of $Q_\sigma(f)$ (see Property 2 in Section 2) is preserved by the multiplication and integration with respect to $v \in \mathbb{R}^d$ of any function monotone function in v .

Although these properties have been shown in many previous analytical works such as those in [25, 49] and others, we will include some constructive proofs of these properties not only for completeness, but also to draw their discrete counterparts when we prove similar properties for the scheme.

We cite several results for existence and uniqueness as well as regularity of initial-boundary value problems associated to the Vlasov-Boltzmann equation (1.1). Y. Guo showed [30, 31] that the non-linear Vlasov Boltzmann Poisson Maxwell system, under spatial periodic boundary conditions and initial data near a global Maxwellian distribution, propagates the regularity of the initial behavior; and further, with R. Strain [47], calculated almost exponential decay rates to such Maxwellian equilibrium.

More recently in case of the initial and boundary value problem, N. Ben-Abdallah and M. Tayeb [2] showed existence and uniqueness of solutions to the linear Vlasov Boltzmann with a continuous in space-time field $E(x, t)$ and non-negative initial and boundary conditions having the same polynomial decay in $L^1 \cup L^\infty$ in one space dimension and higher dimensional phase-space (velocity). This solution preserves the regularity and decay properties of the initial state. While this result uses low regularity of the integrating characteristic field $E(x, t)$ with non vanishing gradients, it is a hope that higher order Sobolev regularity may propagate for more regular fields, as well as more regular initial and boundary conditions satisfying at least polynomial decay,

however such result is still not available. We also mention that in the same manuscript [2] the authors showed the existence of weak solutions to Boltzmann-Poisson system with incoming data with polynomial decay in the case of one phase-space dimension.

We also mention an interesting result of M. Portelheiro [42], who showed that the propagation properties of $v\partial_x + \partial_x V(x)\partial_v + Q$ yields that the closure in $L^1(\mathbb{R}^{2d})$ from C_0^∞ generates a semi-group of contraction in $L^1(\mathbb{R}^{2d})$.

More relevant to our problem of error analysis and long time numerical behavior, is the work of F. Herau [33] who studied the exponential decay properties of the solution of the initial value problem to the stationary solution (1.12) in the case of linear relaxation operator (1.5), by introducing the additional operator

$$\Lambda^2 = -\gamma\partial_v(\partial_v + v) - \gamma\partial_x(\partial_x + \partial_V(x) + 1), \quad (1.14)$$

and showing that problems (1.1) and (1.14) have nice properties in the weighted space

$$B^2 = \{f \in \mathcal{D}' : \frac{f}{\mathcal{M}^{-1/2}} \in L^2(dx, dv)\} \quad (1.15)$$

with the natural norm defined by $\|f\|_{B^2}^2 = \int |f|^2 \mathcal{M}^{-1} dx dv$. Indeed, it is possible to see that the closure from C_0^∞ of the operator $\Lambda^2 - 1$ in the space B^2 is maximal accretive (see [32]) and has 0 as the single eigenvalue associated with the eigenfunction \mathcal{M} and, in addition has a spectral gap $\lambda > 0$ in B^2 . We recall that the spectral gap is defined as the infimum of the spectrum except for the lowest eigenvalue. Example cases when this property holds are either when $\text{Hess}V \geq \lambda I_d$, then λ is the spectral gap; or for $|V'(x)|$ going to infinity in x , then the operator Λ^2 is compact with resolvent in B^2 and so the operator $\Lambda^2 - 1$ has a spectral gap $\lambda > 0$ in B^2 . We refer to [32], [34] and references therein.

In addition it was shown in [33] that in the case of a relaxation operator $Q(f)(v) = \rho(t, x) \mu_\infty(v) - f(x, v)$ the above properties hold in B^2 as well, and the Cauchy problem is well posed and, further, there exists a constant $A > 0$ depending only on the second and third order derivatives of $V(x)$, such that for all L^1 -normalized initial state function $f_0 \in B^2$,

$$\|f(t, x, v) - \mathcal{M}(x, v)\|_{B^2(\mathbb{R}^{2d})} \leq 3 \exp(-\lambda t) \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})}, \quad (1.16)$$

for $f = f(t, x, v)$ the unique solution of equation (1). This is a direct consequence of the decrease of the so-called relative entropy for any $f_0 > 0$. That is

$$0 \leq H(f, \mathcal{M})(t) := \int \int f(t, x, v) \ln \frac{f(t, x, v)}{\mathcal{M}(x, v)} dx dv \leq 3 \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})} \exp(-\lambda t). \quad (1.17)$$

We point out that, if the potential V satisfies the growth condition (1.13), the following functional estimates hold for any function $g \in B^2(\mathbb{R}^{2d})$

$$\|g\|_{L_m^2(\mathbb{R}^{2d})} = \int_{\mathbb{R}_x^d} \int_{\mathbb{R}_v^d} |g|^2 (1 + |x|^2 + |v|^2)^{m/2} dx dv \leq C_V \|g\|_{B^2(\mathbb{R}^{2d})}, \quad (1.18)$$

where $C_V = C_V(K_1, R, m)$ and K_1 and R are constants from (1.13).

Regarding Sobolev regularity results for perturbative states from equilibrium, a broad work of C. Mouhot and N. Newman [38] was also done around the same period

of Herau's and Ben Abdallah-Tayeb's results from [33, 2] on existence and regularity associated with the linear Vlasov Boltzmann equation (1.1). The authors in [38] study the existence, uniqueness regularity and decay rates for a large general class of linear collisional kinetic models in the torus, including in particular the linear collisional integral associated with the linearized Boltzmann equation for hard spheres, the linearized Landau equation with hard and moderately soft potentials and the semi-classical linearized fermionic and bosonic relaxation models. More specifically, they showed explicit coercivity estimates on the associated integro-differential operators for some modified Sobolev norms. They also obtained existence of classical solutions near equilibrium for the full nonlinear models associated with explicit regularity bounds and estimates on the rate of exponential convergence towards equilibrium in the perturbative setting. Their proofs follow from the ideas of coercivity and hypoellipticity as developed in [34, 32, 33, 50] into a characterization of existence, uniqueness, regularity and decay rates for equilibrium perturbative solutions by a study linear transport equations based on a linear energy method. This method is characterized by three properties, namely its mixing in velocity, the splitting of the collisional terms in a coercive form and a regularizing part, and a relaxation term towards an equilibrium (the so called Maxwellian).

More specifically, following [38], the solution to the initial value problem associated to (1.1), with periodic boundary conditions Boltzmann equation, propagates Sobolev regularity and decay estimates. Namely, given the initial state f_0 such that for $\mathcal{M}(x, v)$ the unique stationary state from (1.12)

$$\|\mathcal{M}^{-1/2}(f_0 - \mathcal{M})\|_{H^k} \leq \varepsilon \quad (1.19)$$

for some $k \geq k_0$ and some $0 < \varepsilon \leq \varepsilon_0$ where ε_0 depends explicitly on the collision operator (i.e. the scattering function σ), then there exists a unique global non-negative solution $f = f(t, x, v) \in C([0, \infty), H^k)$ of the initial value problem, such that

$$\|\mathcal{M}^{-1/2}(f(t, \cdot, \cdot) - \mathcal{M})\|_{H^k} \leq C \exp(-\gamma t) \quad \text{for all } t \geq 0 \quad (1.20)$$

with some constants C and $\gamma > 0$. that can be explicitly computed. In particular C is proportional to the initial $\|f_0 - \mathcal{M}\|_{H^k}$ deviation from the equilibrium state \mathcal{M} . The conclusion still holds true when a repulsive self-consistent Poisson potential is added (still with periodic boundary conditions).

Remark: While the cited Sobolev regularity results from [38] required some initial closeness to the equilibrium state \mathcal{M} , the one from [33] does not. It is not known up to think point if the regularity results of [38] will hold for large data. Nevertheless our estimates are for large data, but the computational boundary in phase space is based on the fact that for sufficiently confined potential associated to the field $E(x, t)$ the solution will decay to the unique equilibrium state \mathcal{M} (see subsection immediately below). We also point out that our estimates are done for zero spatial incoming boundary conditions but they are also valid for periodic boundary conditions as well.

1.2. Selection of the computational domain for global in time estimates.

In order to fix the computational domain for our estimates, we proceed as follows. First we notice that the equilibrium state associated to the steady state problem with a linear collisional integral satisfies

$$\|\mathcal{M}\|_{B^2(\Omega)} = \|\mathcal{M}\|_{L^1(\Omega)}^{1/2} \quad (1.21)$$

for any arbitrary set Ω and that, for any arbitrary $\epsilon > 0$, there is an $\Omega_\epsilon \in \mathbb{R}^{2d}$ such that

$$\int_{\mathbb{R}_x^d \times \mathbb{R}_v^d \setminus \Omega_\epsilon} |\mathcal{M}| < \epsilon. \quad (1.22)$$

Essentially the set Ω_ϵ is not “very big” since $f \in B^2$ means that f^2 decays very fast as it is integrable when multiplied by the inverse of the L^1 integrable stationary state $\mathcal{M}(x, v)$. In the particular case where \mathcal{M} is given by (1.12), $\text{diam}(\Omega_\epsilon) \approx -\log \epsilon$ for $\epsilon \ll 1$.

Under the assumption that there is a stationary state $\mathcal{M} \in L^1(\mathbb{R}^{2d})$, from (1.16) for which it is possible to obtain a controlling inequality that yields an stable decay estimate to the stationary state of the form $\|f(x, v, t) - \mathcal{M}\|_{L^2_\mu(\mathbb{R}^{2d})} \leq g(t)\|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})}$ for initial state f_0 and a positive, bounded $g(t)$ such that $\lim_{t \rightarrow \infty} g(t) = 0$. (In fact, due to [33, 38], one may take $g(t) = e^{-\lambda t}$, $\lambda > 0$ in the case of a quadratically confined $V(x)$).

Now we consider the error made by working on the cut-off domain Ω_ϵ . Suppose that the solution $f(x, v, t)$ is uniformly controlled in time and stable with respect to the initial state f_0 for problem (1.1) in all of $\mathbb{R}_+ \times \mathbb{R}^{2d}$. Then, one can estimate the $B^2(\Omega_\epsilon)$ -norm of the solution by

$$\|f\|_{B^2(\Omega_\epsilon)} \leq g(t)\|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})} + \|\mathcal{M}\|_{L^1(\Omega_\epsilon)}^{1/2} \leq \mathbf{K},$$

where the constant \mathbf{K} is uniform in time since $g(t)$ is uniformly bounded in t . Similarly,

$$\|f\|_{B^2(\mathbb{R}^{2d} \setminus \Omega_\epsilon)} \leq g(t)\|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})} + \|\mathcal{M}\|_{L^1(\mathbb{R}^{2d} \setminus \Omega_\epsilon)}^{1/2} = Cg(t) + \epsilon^{1/2},$$

uniformly in time, where $C = \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^{2d})}$, and $\|\mathcal{M}\|_{L^1(\mathbb{R}^{2d} \setminus \Omega_\epsilon)} \leq \epsilon^{1/2}$.

In particular, using that $\lim_{t \rightarrow \infty} g(t) = 0$ there exists a T^* sufficiently large depending on the B^2 -norm distance between the initial and stationary states as well as on the decay rate $g(t)$, such that $C|g(t)| = \mathcal{O}(\epsilon^{1/2})$ for any $t \geq T^*$. Therefore,

$$\|f\|_{B^2(\mathbb{R}^{2d})}(t) = \|f\|_{B^2(\Omega_\epsilon)}(t) + \mathcal{O}(\epsilon^{1/2}), \quad (1.23)$$

or equivalently

$$\|f\|_{B^2(\mathbb{R}^{2d} \setminus \Omega_\epsilon)}(t) = \mathcal{O}(\epsilon^{1/2}), \quad (1.24)$$

uniformly, for any time $t \geq T^*$. In particular, the amount of mass lost by working in the cut-off domain can be controlled by reducing the distance between the initial state and stationary states, or increasing the size of Ω_ϵ . Mass conservation in this domain is also improved if the time decay rate of $g(t)$ is high.

Therefore, it is possible to choose the domain Ω_ϵ big enough, such that it contains almost all of the total mass of the initial datum f_0 and of the stationary states \mathcal{M} from (1.12). Then, at least computationally and well beyond machine accuracy, the solution f of the Cauchy problem Ω_ϵ will take values at the at the boundary of order $\mathcal{O}(\epsilon^{1/2})$ in the B^2 -norm, and so well approximated by be zero, with zero derivatives, and consequently the associated evolution problem will essentially be confined to the domain Ω_ϵ . Thus, choosing the computational domain Ω_ϵ yields that all approximations will have a fixed error $\mathcal{O}(\epsilon^{1/2})$ uniformly in time $t > T^*$, in the $B^2(\Omega_\epsilon)$ -norm.

Furthermore, since the problem is conservative in $L^1(\Omega_\epsilon)$, extending the domain such that the initial data, itself extended with zero values outside Ω_ϵ , is supported at an $\mathcal{O}(1)$ distance from the new boundary, the solution will remain close to the solution of the problem in all space. To the best of our knowledge, there is no available analytical result at the present time to rigourously justify this last statement, which is an assumption for the initial boundary value problem under consideration and the corresponding one in all space.

We remind the reader that the above estimates do not provide a pointwise control of the solution to f outside the domain Ω_ϵ . However we can obtain good estimates in the L^2 -norm. It is important to note that this approach is intended to heuristically justify the selection of the computational domain. However, the calculation of error estimates in the following sections are with respect to the solution of the initial value problem in the bounded domain. Our error estimates are done for zero spatial incoming boundary conditions which are also valid for periodic boundary conditions. We also assume null phase space boundary conditions in the cut-off domain.

2. Properties of the collisional operator. In this section we present some important properties of the linear Boltzmann collision operator that are used our subsequent analysis. As pointed out in Section 1.3, we have to implement numerical calculations on a finite domain in the phase space. The collisional operator under this assumption which will be used in the rest of the paper is

$$Q_\sigma(f)(t, x, v) = \int_{v' \in \mathcal{D}_v^d} (\sigma(x, v, v')f' - \sigma(x, v', v)f)dv', \quad (2.1)$$

The most crucial property of the characterization of the linear collisional operator is related to the scattering function $\sigma(x, v, v')$. Its symmetry and integrability properties determine the conservation and monotonicity properties of the collisional integral, which implies the L^1 -propagation property of the associated Boltzmann transport equation. In particular, this yields a unique positive solution existence result for any transport equation [25]. We include these elementary proofs not only for the sake of completeness, but also because we will extend them to a semi-discrete version of mass conservation of the DG scheme, and in particular to the monotonicity and positivity preservation properties of the first order DG scheme.

PROPERTY 1. (*Symmetry and integrability of the scattering transition probabilities*) *The scattering rate function $\sigma(x, v, v')$, restricted to the computational domain $\Omega_{\mathcal{D}} = \mathcal{D}_x^d \times \mathcal{D}_v^d \subseteq \mathbb{R}_x^d \times \mathbb{R}_v^d$ as defined in section 1.3, is assumed to be positive, to be x -space anisotropic, to satisfy the detailed balanced principle (1.3) and to be integrable in $v \in \mathcal{D}_v^d$. That is,*

$$\begin{aligned} \sigma(x, v, v') &= k(x, v, v') M(v), & k(x, v', v) &= k(x, v, v'), \\ \nu(x, v) &= \int_{\mathcal{D}_{v'}^d} \sigma(x, v', v) dv' < K, & \text{the collision frequency.} \end{aligned}$$

Next, we demonstrate a monotonicity property of any linear collision integral $Q_\sigma(f)$ with $\sigma(x, v, v')$ any (v, v') -space symmetric in $\mathcal{D}_v^d \times \mathcal{D}_{v'}^d$, satisfying *Property 1*. All of the following results are a trivial observation of the fact that solutions to the initial value problem to the linear transport problem produce a monotone mass preserving map. This presentation follows the original work of [25] on relationships between non-expansive and order preserving mappings. This is very natural framework to deal with kinetic collisional transport equations.

PROPERTY 2. (*Monotonicity*) Let $\mathcal{G}(x)$ be any monotone non-decreasing real valued function defined on \mathbb{R} . The following monotonicity formula holds:

$$\int_{\mathcal{D}_v^d} Q_\sigma(f) \mathcal{G}\left(\frac{f}{M}\right) dv \leq 0, \quad \text{for all } x \in \mathcal{D}_x^d. \quad (2.2)$$

Proof. Denote $M = M(v)$ and $M' = M(v')$. By the definition of the linear collision integral and *Property 1* on the positivity and symmetry of the scattering function, we have that

$$\begin{aligned} \int_{\mathcal{D}_v^d} Q_\sigma(f) \mathcal{G}\left(\frac{f}{M}\right) dv &= \int_{\mathcal{D}_v^d \times \mathcal{D}_{v'}^d} (k(x, v, v') M f' - k(x, v', v) M' f) \mathcal{G}\left(\frac{f}{M}\right) dv' dv \\ &= \int_{\mathcal{D}_v^d \times \mathcal{D}_{v'}^d} k(x, v, v') M M' \left(\frac{f'}{M'} - \frac{f}{M} \right) \mathcal{G}\left(\frac{f}{M}\right) dv' dv \\ &= \int_{\mathcal{D}_v^d \times \mathcal{D}_{v'}^d} k(x, v', v) M' M \left(\frac{f}{M} - \frac{f'}{M'} \right) \mathcal{G}\left(\frac{f'}{M'}\right) dv' dv \\ &= \frac{1}{2} \int_{\mathcal{D}_v^d \times \mathcal{D}_{v'}^d} k(x, v, v') M M' \left(\frac{f'}{M'} - \frac{f}{M} \right) \left(\mathcal{G}\left(\frac{f}{M}\right) - \mathcal{G}\left(\frac{f'}{M'}\right) \right) dv' dv \\ &\leq 0. \end{aligned}$$

by the symmetry of $k(x, v, v')$ and by the monotonicity of \mathcal{G} . \square

An immediate consequence of the monotonicity *Property 2* is conservation of mass, which is obtained by taking the function $\mathcal{G}(x) = 1$, and observing that the above proof yields an identity.

PROPERTY 3. (*Conservation*) The collisional operator is mass conservative:

$$\int_{\mathcal{D}_v^d} Q_\sigma(f) dv = 0.$$

The next property is also a direct application of the monotonicity *Property 2*, and in fact yields a classical type of so-called entropy or energy estimate for first order transport equations when multiplied by monotone non-decreasing functions with convex antiderivatives. This is indeed also a consequence the non-expansive measure preserving nature of the collisional integral as also shown in [25].

PROPERTY 4. (*L^1 Contraction*) Assuming zero boundary conditions, the collisional operator is L^1 -contractive in the sense that

$$\frac{d}{dt} \|f\|_{L^1(\mathcal{D}_v^d \times \mathcal{D}_x^d)} \leq 0.$$

Proof. Note that $|f|_t = f_t \operatorname{sgn}(f) = f_t \operatorname{sgn}\left(\frac{f}{M}\right)$, and $\nabla|f| = \nabla f \operatorname{sgn}(f) = \nabla f \operatorname{sgn}\left(\frac{f}{M}\right)$. Multiply equation (1.9) by $\mathcal{G}(f) = \operatorname{sgn}\left(\frac{f}{M}\right)$, which is a monotone increasing function of its argument, and integrate over $\mathbb{R}_x^d \times \mathbb{R}_v^d$ to obtain

$$\begin{aligned} \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \frac{\partial f}{\partial t} \operatorname{sgn}\left(\frac{f}{M}\right) dv dx + \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \alpha \cdot \nabla f \operatorname{sgn}\left(\frac{f}{M}\right) dv dx \\ = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} Q_\sigma(f) \operatorname{sgn}\left(\frac{f}{M}\right) dv dx \leq 0, \quad (2.3) \end{aligned}$$

by *Property 2*, and

$$\int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \alpha \cdot \nabla f \operatorname{sgn}\left(\frac{f}{M}\right) dv dx = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \alpha \cdot \nabla |f| dv dx = 0, \quad (2.4)$$

and consequently,

$$\frac{d}{dt} \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} |f| dv dx = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \frac{\partial f}{\partial t} \operatorname{sgn}\left(\frac{f}{M}\right) dv dx \leq 0 \quad (2.5)$$

thus concluding the proof. \square

A natural corollary of the monotonicity *Property 2* and mass preserving *Property 3* is the positivity of the solution to the initial value problem.

PROPERTY 5. (*Positivity of the solution to the initial value problem*) Assuming zero boundary conditions, the solution to the initial value problem (1.1) is positive for all times if the initial probability $f_0 = f(0, x, v)$ is positive.

Proof. The proof is very similar to the one for *Property 4* above. Here choose $\mathcal{G} = \frac{1}{2}(\operatorname{sgn}(\frac{f}{M}) - 1) = \frac{1}{2}(\operatorname{sgn}(f) - 1)$ as a test function for the linear Boltzmann equation. Since the negative part of f , defined as $f^- = \max\{0, -f\}$ is the antiderivative of \mathcal{G} , as in (2.3) and (2.4),

$$\frac{d}{dt} \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} f^- dv dx = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} Q_\sigma(f) \frac{\operatorname{sgn}(\frac{f}{M}) - 1}{2} dv dx \leq 0. \quad (2.6)$$

Thus one obtains that

$$\int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} f^- dv dx \leq \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} f^-(0, x, v) dv dx = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \max\{0, -f_0\} dv dx = 0, \quad (2.7)$$

since f_0 is always taken positive. Since the integrand f^- is a non-negative function whose integral, computed in (2.7), was shown non-positive, then it must be $f^-(t, x, v) = 0$ for all $t \geq 0$. Consequently, $f(t, x, v)$ is positive for all $t \geq 0$. \square

Remark: In fact, the positivity and mass conservation properties immediately yield conservation of the L^1 norm, and thus L^1 stability as well. Although the above proof may be viewed as redundant, we have included it for completeness for our numerical scheme properties. However, the resulting scheme is stable for any order of the approximating polynomial space. Obviously all these estimates have an error $O(\epsilon)$ with respect to the continuum solution in the whole space.

The next lemma will provide $L^2(\Omega_{\mathcal{D}})$ control estimates for the collision operator under suitable assumptions on the growth of the scattering rate function $\sigma(x, v, v')$ restricted to $\mathcal{D}_x^d \times \mathcal{D}_v^d \times \mathcal{D}_{v'}^d$. It is important for the stability and error estimates for the DG scheme.

LEMMA 6. ($L^2(\Omega_{\mathcal{D}})$ control of the collisional integral) Assume the scattering rate function satisfies (2.2) and

$$0 \leq \sigma(x, v, \cdot) \leq C_1 + C_2(|x| + |v|). \quad (2.8)$$

Then for f and $g \in L^2(\Omega_{\mathcal{D}})$, the following estimate hold

$$\int_{\Omega_{\mathcal{D}}} Q_\sigma(f) g dx dv \leq C_\sigma \operatorname{diam}(\Omega_{\mathcal{D}}) \operatorname{diam}(\mathcal{D}_v^d) \|f\|_{L^2(\Omega_{\mathcal{D}})} \|g\|_{L^2(\Omega_{\mathcal{D}})}. \quad (2.9)$$

These estimates are time dependent and their parameters given by $C_\sigma = 2 \max\{C_1, 2C_2\}$ depend on σ .

Remark: We note that the growth condition (2.8) can be taken as the same growth of the advection vector $\alpha = (v, \nabla V)$. That means the growth of the scattering function σ as a function of x , or its corresponding collision frequency function ν , can be of the same order as the growth of the gradient of the potential. In our particular problem, due to condition (1.13), we assume *at most* linear growth in x .

Proof.

$$\begin{aligned} & \int_{(x,v) \in \Omega_{\mathcal{D}}} Q_\sigma(f) g \, dv dx \\ &= \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \int_{\mathcal{D}_v^d} (\sigma(x, v, v') f(x, v') - \sigma(x, v', v) f(x, v)) g(x, v) \, dv' dv dx \\ &= A^1 - A^2, \end{aligned}$$

where

$$A^1 = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \int_{\mathcal{D}_v^d} \sigma(x, v, v') f(x, v') g(x, v) \, dv' dv dx,$$

and

$$A^2 = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \int_{\mathcal{D}_v^d} \sigma(x, v', v) f(x, v) g(x, v) \, dv' dv dx.$$

Since $\sigma(x, v, \cdot) \leq C_1 + C_2(|x| + |v|)$, we have

$$\begin{aligned} |A^1| &\leq \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \int_{\mathcal{D}_v^d} (C_1 + C_2(|x| + |v|)) |f(x, v')| |g(x, v)| \, dv' dv dx \\ &\leq \frac{C_\sigma}{2} \text{diam}(\Omega_{\mathcal{D}}) \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} \int_{\mathcal{D}_v^d} |f(x, v')| |g(x, v)| \, dv' dv dx \\ &\leq \frac{C_\sigma}{2} \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) \int_{\mathcal{D}_x^d} \|f(x, \cdot)\|_{L^2(\mathcal{D}_v^d)} \cdot \|g(x, \cdot)\|_{L^2(\mathcal{D}_v^d)} \, dx \\ &\leq \frac{C_\sigma}{2} \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) \|f\|_{L^2(\Omega_{\mathcal{D}})} \cdot \|g\|_{L^2(\Omega_{\mathcal{D}})} \end{aligned}$$

By similar arguments,

$$|A^2| \leq \frac{C_\sigma}{2} \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) \|f\|_{L^2(\Omega_{\mathcal{D}})} \cdot \|g\|_{L^2(\Omega_{\mathcal{D}})},$$

so

$$\int_{\Omega_{\mathcal{D}}} Q_\sigma(f) g \, dx dv \leq C_\sigma \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) \|f\|_{L^2(\Omega_{\mathcal{D}})} \cdot \|g\|_{L^2(\Omega_{\mathcal{D}})}.$$

□

3. Formulation of the DG Scheme. In this section, and without loss of generality, we assume $\mathcal{D}_v^d = [-V_i, V_i]^d$ and $\mathcal{D}_x^d = [0, L_i]^d$, for $0 < V_i, L_i < \infty$, $i = 1 \dots d$.

Let $\{\mathcal{T}_h\}$ denote a family of a non degenerate finite element subdivisions of $\Omega_{\mathcal{D}}$ partitioned into open disjoint elements K with exterior boundary $\partial\Omega_{\mathcal{D}}$. We denote h_K to be the diameter of element K , ρ_K to be the diameter of the biggest sphere included in K , we impose the classical assumption of shape regularity [15], $h_K/\rho_K \leq \sigma_0$, and let $\mathbf{h} = \sup_K h_K$. Denote the set of all element edges associated to this mesh as $\mathfrak{e}_h = \cup_{K \in \mathcal{T}_h} \partial K$. \mathfrak{e}_h is defined to allow redundancy. For example, if an edge $e \in \partial K_1$ and $e \in \partial K_2$, it will appear twice. This notation allows for the hanging nodes in the mesh. Denote edges that belongs to $\Omega_{\mathcal{D}}$ as

$$\begin{aligned} F^{0-} (F^{0+}) & \quad \text{the set of faces located on } \partial\Omega_{\mathcal{D}} \text{ such that } x = 0, v < 0 (v > 0), \\ F^{L_i-} (F^{L_i+}) & \quad \text{the set of faces located on } \partial\Omega_{\mathcal{D}} \text{ such that } x = L_i, v < 0 (v > 0), \\ F^{-V_i} (F^{+V_i}) & \quad \text{the set of faces located on } \partial\Omega_{\mathcal{D}} \text{ such that } v = -V_i (v = +V_i). \end{aligned}$$

Define the inflow face as $\Gamma^- = \cup_i \{F^{0+} \cup F^{L_i-}\}$, and suppose $f_h = f^{in}$ on Γ^- . Define $\Gamma^0 = \cup_i \{F^{-V_i} \cup F^{+V_i}\}$. Since f is a probability distribution, it is reasonable to enforce $f_h = 0$ on Γ^0 when V_i is large enough.

The finite element space is defined as

$$V_h^k = \{\phi_h \in L^2(\Omega_{\mathcal{D}}) : \forall K \in \mathcal{T}_h(\Omega_{\mathcal{D}}), \phi_h|_K \in P^k(K)\} \quad (3.1)$$

where $P^k(K)$ is the set of polynomials of total degree at most k on the simplex K .

3.1. The semi-discrete DG scheme. The semi-discrete discontinuous Galerkin discretization of (1.9) is given as follows: seek $f_h = f_h(t, x, v) \in \mathbb{R}_t^+ \times V_h^k$ such that, the below equality holds for all for all test function $w_h \in V_h^k$,

$$(\partial_t f_h, w_h)_{\mathcal{T}_h} + \mathcal{A}(f_h, w_h) = \mathcal{L}(w_h), \quad (3.2)$$

where

$$\mathcal{A}(f_h, w_h) \equiv -(f_h, \alpha \cdot \nabla w_h)_{\mathcal{T}_h} + \langle \hat{f}_h, w_h \alpha \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h \setminus \Gamma^- \setminus \Gamma^0} - (Q_\sigma(f_h), w_h)_{\mathcal{T}_h} \quad (3.3)$$

and

$$\mathcal{L}(w_h) \equiv -\langle f^{in}, w_h \alpha \cdot \mathbf{n} \rangle_{\Gamma^-} \quad (3.4)$$

In the above equalities, we are using the following notations

$$(\zeta, w)_{\mathcal{T}_h} \equiv \sum_{K \in \mathcal{T}_h} \int_K \zeta w dx dv, \quad \langle \zeta, w \alpha \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} \equiv \sum_{K \in \mathcal{T}_h} \int_{\partial K} \zeta(\gamma) w(\gamma) \alpha(\gamma) \cdot \mathbf{n} d\gamma,$$

and \mathbf{n} denotes the outward unit normal vector to ∂K . For any edge $e \in \mathfrak{e}_h$ that is not associated with a particular K , \mathbf{n} can be defined as the unit normal in either of the possible direction. \hat{f}_h is the monotone numerical flux, which is chosen as the upwind flux in our case.

$$\hat{f}_h = f_h^-,$$

where $f_h^-(z) = \lim_{\delta \downarrow 0} f_h(z - \delta \alpha(z))$, and $\alpha(\cdot)$ is defined in (1.10).

The standard L^2 projection for any function u is a function $\mathbb{P}u \in V_h^k$, such that

$$(\mathbb{P}u - u, w_h)_{\mathcal{T}_h} = 0,$$

for all $w_h \in V_h^k$. The initial condition is defined through the L^2 projection $f_h(0, x, v) = \mathbb{P}f(0, x, v)$.

In practice, the face and volume integrals in $\mathcal{A}(f_h, \phi_h)$, $\mathcal{L}(\phi_h)$ need to be evaluated by certain numerical quadratures. One can choose to implement a very high order quadrature on those terms to guarantee the quadrature error is on the order of machine precision. Furthermore, since the scheme is linear on f_h , those integrations only need to be performed on the basis functions and can be stored before the time evolution starts. On the other hand, we can also use quadrature formulas with fixed number of points. In [18], it was proven that a quadrature formula for face integrals that is exact for polynomials of degree $(2k + 1)$ and element integrals that is exact for degree $2k$ polynomials can guarantee an error in the $\|\cdot\|_\infty$ norm of order \mathbf{h}^{k+1} . This error is on the order of the distance of f to the finite element space V_h^k and will not deteriorate the quality of numerical solution. For simplicity of discussion, from this point on, we shall assume the integrals are evaluated exactly for the semi-discrete scheme.

3.2. Time discretization. We use total variation diminishing (TVD) high-order Runge-Kutta methods [46] to solve the method of lines ODE resulting from the semi-discrete scheme,

$$(f_h)_t = R(f_h). \quad (3.5)$$

Those time stepping methods are convex combinations of the Euler forward time discretization. The commonly used third-order TVD Runge-Kutta method is given by

$$\begin{aligned} f_h^{(1)} &= f_h^n + \Delta t R(f_h^n) \\ f_h^{(2)} &= \frac{3}{4} f_h^n + \frac{1}{4} f_h^{(1)} + \frac{1}{4} \Delta t R(f_h^{(1)}) \\ f_h^{n+1} &= \frac{1}{3} f_h^n + \frac{2}{3} f_h^{(2)} + \frac{2}{3} \Delta t R(f_h^{(2)}) \end{aligned} \quad (3.6)$$

Detailed description of the TVD Runge-Kutta method can be found in [46], see also [28], and [29] for strong-stability-preserving method.

3.3. The limiter and positivity-preserving DG scheme. The scheme in Section 3.1 conserves mass. Furthermore, it is monotone and thus positivity-preserving for piecewise constant approximations under certain CFL conditions (cf Section 4). However, it is highly nontrivial how to preserve positivity of the solution for high order discretizations. In [52], a limiter that preserves the maximum principle of conservation laws is proposed. This limiter is uniformly high order accurate. Moreover, it does not change the cell average of the numerical solution. We present the limiter below and defer the detailed analysis of the scheme to Section 5.

In each of the forward Euler step of the time discretization, the following procedures are performed:

- On each simplex K , evaluate $T_K = \min_{(x,v) \in S_K} f_h(x, v)$, where S_K denotes the set of certain points on K that will be described in Section 5.

- Compute $\tilde{f}_h(x, v) = \theta(f_h(x, v) - (\overline{f_h})_K) + (\overline{f_h})_K$, where $(\overline{f_h})_K$ is the cell average of f_h on K , and $\theta = \min\{1, |(\overline{f_h})_K|/|T_K - (\overline{f_h})_K|\}$. This limiter has the effect of keeping the cell average and “squeeze” the function to be positive at points in S_K .
- Use \tilde{f}_h instead of f_h to compute the Euler forward step. A suitable CFL condition will be derived in Section 5.

In [52], it was shown that this limiter is uniformly $(k+1)$ -th order accurate. We defer relevant analysis of the scheme to Section 5.

4. Analysis of the semi-discrete DG scheme. In this section, we will analyze the properties of the semi-discrete DG scheme (3.2). For simplicity of discussion, we assume periodic boundary conditions on f in all directions. The discussion also implies the case of $f^{in} = 0$. The general case will involve terms like $\|f^{in}\|_{L^2(\Gamma^-)}$ in the stability proof and will destroy mass conservation. The proof is rather technical and we do not pursue it in this paper.

If periodic conditions are imposed, then the term $\mathcal{A}(f_h, w_h)$ and $\mathcal{L}(w_h)$ are simplified as

$$\mathcal{A}(f_h, w_h) = -(f_h, \alpha \cdot \nabla w_h)_{\mathcal{T}_h} + \langle f_h^-, w_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), w_h)_{\mathcal{T}_h},$$

and

$$\mathcal{L}(w_h) = 0.$$

4.1. Mass conservation. The following theorem implies the mass conservation property of the semi-discrete DG scheme.

THEOREM 7. *The semi-discrete DG scheme conserves mass, i.e.*

$$\frac{d}{dt} \int_{\Omega_D} f_h dv dx = 0. \quad (4.1)$$

Proof. Let $w_h = 1$ in (3.2). (4.1) follows from Property 3 and $\langle f_h^-, 1 \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = 0$. \square

4.2. L^2 stability and error estimate for arbitrary order. Here we prove the L^2 stability and error estimate for (3.2) with arbitrary order of discretization. The proof use heavily the properties of the operator $Q_\sigma(\cdot)$ derived in Section 2.

4.2.1. The case when $M(v) = \text{constant}$. We first consider the simplified case of $M(v) = \text{constant}$ to illustrate the ideas of the proof. The general case will be treated in Section 4.2.2. When $M(v) = \text{constant}$, $\sigma(x, v, v')$ is symmetric about v and v' and Property 2 is reduced to

$$\int_{\mathcal{D}_v^d} Q_\sigma(f) \mathcal{G}(f) dv \leq 0, \quad \text{for all } \mathcal{G} \text{ that is a monotone non-decreasing function.}$$

THEOREM 8. (*L^2 stability when $M(v) = \text{constant}$*) *Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9) with $M(v) = \text{constant}$, we have,*

$$\|f_h(t)\|_{L^2(\Omega_D)} \leq \|f_h(0)\|_{L^2(\Omega_D)}. \quad (4.2)$$

Proof. Let the test function $w_h = f_h$ in (3.2), then

$$(\partial_t f_h, f_h)_{\mathcal{T}_h} + \mathcal{A}(f_h, f_h) = 0,$$

where

$$\begin{aligned} \mathcal{A}(f_h, f_h) &= -(f_h, \alpha \cdot \nabla f_h)_{\mathcal{T}_h} + \langle f_h^-, f_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), f_h)_{\mathcal{T}_h} \\ &= -\frac{1}{2} \langle f_h, f_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} + \langle f_h^-, f_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), f_h)_{\mathcal{T}_h} \\ &= \frac{1}{4} \langle [f_h], [f_h] | \alpha \cdot \mathbf{n} \rangle_{\mathfrak{e}_h} - (Q_\sigma(f_h), f_h)_{\mathcal{T}_h} \\ &= \frac{1}{4} \| [f_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 - (Q_\sigma(f_h), f_h)_{\mathcal{T}_h}, \end{aligned}$$

where the notation $\langle \zeta, w \rangle_{\mathfrak{e}_h} = \sum_{e \in \mathfrak{e}_h} \langle \zeta, w \rangle_e$, and $[f_h] = f_h^+ - f_h^-$ denotes the jump. Note that in the above equality, the factor is $\frac{1}{4}$ not $\frac{1}{2}$ due to the redundancy in the notation \mathfrak{e}_h . It holds true for the case of α as a variable, and allow $\alpha \cdot \mathbf{n}$ to change sign within a single edge. Because of Property 3 of Q_σ , $(Q_\sigma(f_h), f_h)_{\mathcal{T}_h} \leq 0$. Hence

$$\mathcal{A}(f_h, f_h) \geq 0,$$

and we are done. \square

In order to prove L^2 error estimates, we follow the classic work of [35] for constant coefficient conservation laws. The main difficulty is that we are treating a variable coefficient equation with a collisional integral. The averaging technique in the proof of [17] will be adopted to take care of the variable coefficient α , while the collision term will be bounded using Lemma 7. We remark that the accuracy order $(k + \frac{1}{2})$ obtained below is optimal under general meshes [41]. The error estimates could be improved to $(k + 1)$ -th order, if special restrictions on the mesh and special projections are used, see [45, 19, 16, 17] for related discussions, and we will not pursue it in this paper.

THEOREM 9. (*L^2 error estimate when $M(v) = \text{constant}$*) Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9) with $M(v) = \text{constant}$, we have

$$\|f_h(t, \cdot, \cdot) - f(t, \cdot, \cdot)\|_{L^2(\Omega_D)} \leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_D))}, \quad (4.3)$$

where $C = C(\text{diam}(\Omega_D), \|\alpha\|_{\mathbf{W}^{1,\infty}(\Omega_D)})$ and C does not depend on h or t .

Proof. Since the exact solution f also satisfies the (3.2), we have

$$(\partial_t f, w_h)_{\mathcal{T}_h} + \mathcal{A}(f, w_h) = \mathcal{L}(w_h),$$

for any test function $w_h \in V_h^k$. If we define the error as $\mathbb{E} = f - f_h$, then

$$(\partial_t \mathbb{E}, w_h)_{\mathcal{T}_h} + \mathcal{A}(\mathbb{E}, w_h) = 0.$$

We use the L^2 projection \mathbb{P} to decompose \mathbb{E} into two parts, namely $\mathbb{E} = \mathcal{E} + E_h$, where $\mathcal{E} = f - \mathbb{P}f$ and $E_h = \mathbb{P}f - f_h$. Clearly, $E_h \in V_h^k$, thus

$$(\partial_t \mathbb{E}, E_h)_{\mathcal{T}_h} + \mathcal{A}(\mathbb{E}, E_h) = 0,$$

which implies

$$(\partial_t E_h, E_h)_{\mathcal{T}_h} + \mathcal{A}(E_h, E_h) = -(\partial_t \mathcal{E}, E_h)_{\mathcal{T}_h} - \mathcal{A}(\mathcal{E}, E_h). \quad (4.4)$$

Following the definition of L^2 projection, we have $(\mathcal{E}, w_h)_{\mathcal{T}_h} = 0$ for any $w_h \in V_h^k$. Hence $(\partial_t \mathcal{E}, w_h)_{\mathcal{T}_h} = 0$ for any $w_h \in V_h^k$, and

$$(\partial_t \mathcal{E}, E_h)_{\mathcal{T}_h} = 0. \quad (4.5)$$

Similar to the stability proof

$$\mathcal{A}(E_h, E_h) = \frac{1}{4} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 - (Q_\sigma(E_h), E_h)_{\mathcal{T}_h}. \quad (4.6)$$

Plug (4.5) and (4.6) into (4.4), we get

$$(\partial_t E_h, E_h)_{\mathcal{T}_h} + \frac{1}{4} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 = (Q_\sigma(E_h), E_h)_{\mathcal{T}_h} - \mathcal{A}(\mathcal{E}, E_h) \leq -\mathcal{A}(\mathcal{E}, E_h).$$

Next we will try to bound the right hand side of the above inequality. Since

$$-\mathcal{A}(\mathcal{E}, E_h) = (\mathcal{E}, \alpha \cdot \nabla E_h)_{\mathcal{T}_h} - \langle \mathcal{E}^-, E_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} + (Q_\sigma(\mathcal{E}), E_h)_{\mathcal{T}_h}.$$

If we define

$$\begin{aligned} T_1 &= (\mathcal{E}, \alpha \cdot \nabla E_h)_{\mathcal{T}_h}, \\ T_2 &= \langle \mathcal{E}^-, E_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}, \\ T_3 &= (Q_\sigma(\mathcal{E}), E_h)_{\mathcal{T}_h}, \end{aligned}$$

then

$$-\mathcal{A}(\mathcal{E}, E_h) \leq |T_1| + |T_2| + |T_3|.$$

It remains to estimate the terms T_i , $i = 1, 2, 3$.

Estimate of T_1

For term T_1 , if α is a constant, then $T_1 = 0$ since $\alpha \cdot \nabla E_h \in V_h^k$. However, this is not true when α depends on x and v . In this case, similar to [17], we define an average of α on each simplex K as a constant vector α^0 , such that

$$\langle (\alpha - \alpha^0) \cdot \mathbf{n}, 1 \rangle_{\partial K} = 0.$$

Then it follows that

$$T_1 = (\mathcal{E}, (\alpha - \alpha^0) \cdot \nabla E_h)_{\mathcal{T}_h},$$

which implies

$$|T_1| \leq \sum_{K \in \mathcal{T}_h} \|\mathcal{E}\|_{L^2(K)} \|\alpha - \alpha^0\|_{L^\infty(K)} \|\nabla E_h\|_{L^2(K)},$$

By the inverse inequality, there exists a constant C_K such that

$$\|\nabla E_h\|_{L^2(K)} \leq C_K \|E_h\|_{L^2(K)} / h_K.$$

Hence,

$$\begin{aligned}
|T_1| &\leq \sum_{K \in \mathcal{T}_h} C_K \|\mathcal{E}\|_{L^2(K)} \|E_h\|_{L^2(K)} \{ \|\alpha - \alpha^0\|_{L^\infty(K)} / h_K \} \\
&\leq C \max_{K \in \mathcal{T}_h} \{ \|\alpha - \alpha^0\|_{L^\infty(K)} / h_K \} \sum_{K \in \mathcal{T}_h} \|\mathcal{E}\|_{L^2(K)} \|E_h\|_{L^2(K)}, \\
&\leq C \max_{K \in \mathcal{T}_h} \{ \|\alpha - \alpha^0\|_{L^\infty(K)} / h_K \} \|\mathcal{E}\|_{L^2(\mathcal{T}_h)} \|E_h\|_{L^2(\mathcal{T}_h)}, \\
&\leq C |\alpha|_{\mathbf{W}^{1,\infty}(\mathcal{T}_h)} \|\mathcal{E}\|_{L^2(\mathcal{T}_h)} \|E_h\|_{L^2(\mathcal{T}_h)},
\end{aligned}$$

where the last inequality holds similar to [17]. Now, recall a standard estimate for the L^2 projection,

$$\|\mathcal{E}\|_{L^2(K)} + h_K^{\frac{1}{2}} \|\mathcal{E}\|_{L^2(\partial K)} + h_K \|\nabla \mathcal{E}\|_{L^2(K)} \leq C h_K^{k+1} |f|_{H^{k+1}(K)}.$$

Finally we have

$$|T_1| \leq C |\alpha|_{\mathbf{W}^{1,\infty}(\mathcal{T}_h)} |f|_{H^{k+1}(\mathcal{T}_h)} h^{k+1} \|E_h\|_{L^2(\mathcal{T}_h)}.$$

Estimate of T_2

$$\begin{aligned}
|T_2| &= |\langle \mathcal{E}^-, E_h \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}| = \left| \frac{1}{2} \langle \mathcal{E}^-, (E_h^- - E_h^+) \alpha \cdot \mathbf{n} \rangle \right| \\
&\leq \frac{1}{8} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 + \frac{1}{2} \| \mathcal{E}^- \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 \\
&\leq \frac{1}{8} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 + C \|\alpha\|_{L^\infty(\mathcal{T}_h)} \|\mathcal{E}\|_{L^2(\mathfrak{e}_h)}^2 \\
&\leq \frac{1}{8} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 + C h^{2k+1} |f|_{H^{k+1}(\mathcal{T}_h)}^2 \|\alpha\|_{L^\infty(\mathcal{T}_h)}
\end{aligned}$$

Estimate of T_3

It follows from Lemma 6 that

$$\begin{aligned}
T_3 &= (Q_\sigma(\mathcal{E}), E_h)_{\mathcal{T}_h} = \int_{\mathcal{D}_x^d} \int_{\mathcal{D}_v^d} Q_\sigma(\mathcal{E}) E_h dv dx \\
&\leq C \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) \|\mathcal{E}\|_{L^2(\Omega_{\mathcal{D}})} \|E_h\|_{L^2(\Omega_{\mathcal{D}})} \\
&\leq C \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d) |f|_{H^{k+1}(\mathcal{T}_h)} h^{k+1} \|E_h\|_{L^2(\mathcal{T}_h)}.
\end{aligned}$$

Combining the estimates of T_1 , T_2 and T_3 , we have

$$\begin{aligned}
&(\partial_t E_h, E_h)_{\mathcal{T}_h} + \frac{1}{4} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 \\
&\leq \frac{1}{8} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\mathfrak{e}_h)}^2 + C \|\alpha\|_{L^\infty(\mathcal{T}_h)} h^{2k+1} |f|_{H^{k+1}(\mathcal{T}_h)}^2 \\
&\quad + C(|\alpha|_{\mathbf{W}^{1,\infty}(\mathcal{T}_h)} + \text{diam}(\Omega_{\mathcal{D}}) \text{diam}(\mathcal{D}_v^d)) |f|_{H^{k+1}(\mathcal{T}_h)} h^{k+1} \|E_h\|_{L^2(\mathcal{T}_h)}.
\end{aligned}$$

In the case of quadratically confined electrostatic potentials, $\|\alpha\|_{\mathbf{W}^{1,\infty}(\Omega_{\mathcal{D}})} \leq C$, this implies

$$(\partial_t E_h, E_h)_{\mathcal{T}_h} \leq C |f|_{H^{k+1}(\Omega_{\mathcal{D}})}^2 h^{2k+1} + C |f|_{H^{k+1}(\Omega_{\mathcal{D}})} h^{k+1} \|E_h\|_{L^2(\mathcal{T}_h)},$$

where $C = C(\|\alpha\|_{\mathbf{W}^{1,\infty}(\Omega_{\mathcal{D}})})$. Hence,

$$\frac{d}{dt} \|E_h\|_{L^2(\Omega_{\mathcal{D}})}^2 \leq C h^{2k+1} \|f\|_{H^{k+1}(\Omega_{\mathcal{D}})}^2 + C h \|E_h\|_{L^2(\Omega_{\mathcal{D}})}^2.$$

Since $E_h = 0$ at $t = 0$, we have

$$\|E_h\|_{L^2(\Omega_{\mathcal{D}})} \leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_{\mathcal{D}}))},$$

so

$$\|\mathbb{E}\|_{L^2(\Omega_{\mathcal{D}})} \leq C (\sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_{\mathcal{D}}))} + h^{k+1}).$$

We remark that, if $\|f(t=0)\|_{H^{k+1}(\Omega_{\mathcal{D}})} \leq C$, then $\|f(t)\|_{H^{k+1}(\Omega_{\mathcal{D}})} \leq C$, by the regularity propagation property (1.20) (see [38] for detailed discussions.) For the practically relevant case of $t \geq Ch$, i.e., after a few time steps, we obtain the L^2 error estimates of the DG solution as

$$\|\mathbb{E}\|_{L^2(\Omega_{\mathcal{D}})} \leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_{\mathcal{D}}))}.$$

□

4.2.2. The general case when $M(v)$ is not a constant. In the case of general $M(v)$, $\int_{\mathcal{D}_v^d} Q_\sigma(f) f dv \leq 0$ is no longer true. This will introduce some extra terms in the proof. We can see from below the L^2 norm of the f_h is no longer strictly decaying.

THEOREM 10. (*L^2 stability for general $M(v)$*) Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9), we have,

$$\|f_h(t)\|_{L^2(\Omega_{\mathcal{D}})} \leq \exp(Ct) \|f_h(0)\|_{L^2(\Omega_{\mathcal{D}})}, \quad (4.7)$$

where $C = C(\text{diam}(\Omega_{\mathcal{D}}))$.

Proof. The proof follows similar to those in Theorem 8. Let the test function $w_h = f_h$ in (3.2), then

$$(\partial_t f_h, f_h)_{\mathcal{T}_h} + \mathcal{A}(f_h, f_h) = 0,$$

where

$$\mathcal{A}(f_h, f_h) = \frac{1}{4} \| [f_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\epsilon_h)}^2 - (Q_\sigma(f_h), f_h)_{\mathcal{T}_h}.$$

Now $(Q_\sigma(f_h), f_h)_{\mathcal{T}_h}$ is not necessarily non-negative. Here we use Lemma 7 again, and we have

$$\frac{1}{2} \frac{d}{dt} \|f_h\|_{L^2(\Omega_{\mathcal{D}})}^2 = (\partial_t f_h, f_h)_{\mathcal{T}_h} \leq (Q_\sigma(f_h), f_h)_{\mathcal{T}_h} \leq C \|f_h\|_{L^2(\Omega_{\mathcal{D}})}^2,$$

where $C = C(\text{diam}(\Omega_{\mathcal{D}}))$. This implies

$$\|f_h(t)\|_{L^2(\Omega_{\mathcal{D}})} \leq \exp(Ct) \|f_h(0)\|_{L^2(\Omega_{\mathcal{D}})},$$

and we are done.

□

Remark: We can design a semi-discrete DG scheme for (1.9) to preserve the strict decay of the B^2 norm for f_h , in fact if we seek $f_h = f_h(t, x, v) \in \mathbb{R}_t^+ \times V_h^k$ such that, the below equality holds for all for all test function $w_h \in V_h^k$,

$$(\partial_t f_h, \frac{w_h}{\mathcal{M}(x, v)})_{\mathcal{T}_h} + \mathcal{A}(f_h, \frac{w_h}{\mathcal{M}(x, v)}) = \mathcal{L}(\frac{w_h}{\mathcal{M}(x, v)}),$$

then, by the same argument as in the proof of Theorem 9 and the fact that $\alpha \cdot \nabla_{x,v} \mathcal{M}(x, v) = 0$, we can easily deduce $\|f_h(t)\|_{B^2(\Omega_D)} \leq \|f_h(0)\|_{B^2(\Omega_D)}$. However, we choose not to use this scheme because \mathcal{M} will be near zero in a large portion of the domain, which will cause difficulty in numerical implementation.

THEOREM 11. (*L^2 error estimate for general $M(v)$*) Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9), we have

$$\|f_h(t, \cdot, \cdot) - f(t, \cdot, \cdot)\|_{L^2(\Omega_D)} \leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_D))}, \quad (4.8)$$

where $C = C(\text{diam}(\Omega_D), \|\alpha\|_{\mathbf{W}^{1,\infty}(\Omega_D)})$ and C does not depend on h or t .

Proof. The proof is very similar to those in Theorem 9. Using the same set of notation, we have

$$(\partial_t E_h, E_h)_{\mathcal{T}_h} + \frac{1}{4} \| [E_h] \sqrt{|\alpha \cdot \mathbf{n}|} \|_{L^2(\epsilon_h)}^2 = (Q_\sigma(E_h), E_h)_{\mathcal{T}_h} - \mathcal{A}(\mathcal{E}, E_h),$$

except the fact that $(Q_\sigma(E_h), E_h)_{\mathcal{T}_h}$ is not necessarily negative. Following the same lines of proof to bound $-\mathcal{A}(\mathcal{E}, E_h)$, we obtain

$$\begin{aligned} \frac{d}{dt} \|E_h\|_{L^2(\Omega_D)}^2 &\leq C h^{2k+1} \|f\|_{H^{k+1}(\Omega_D)}^2 + C h \|E_h\|_{L^2(\Omega_D)}^2 + (Q_\sigma(E_h), E_h)_{\mathcal{T}_h} \\ &\leq C h^{2k+1} \|f\|_{H^{k+1}(\Omega_D)}^2 + C h \|E_h\|_{L^2(\Omega_D)}^2 + C \|E_h\|_{L^2(\Omega_D)}^2 \\ &\leq C h^{2k+1} \|f\|_{H^{k+1}(\Omega_D)}^2 + C \|E_h\|_{L^2(\Omega_D)}^2. \end{aligned}$$

Hence (4.8) follows using $E_h(t=0) = 0$. \square

The next theorem concerns the decay of numerical solution f_h towards equilibrium. It's a direct consequence of the error estimates above.

THEOREM 12. (*L^2 decay of numerical solution towards equilibrium*) Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9), we have

$$\|f_h(t, \cdot, \cdot) - \mathcal{M}\|_{L^2(\Omega_D)} \leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_D))} + 3e^{-\lambda t} \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^d)}, \quad (4.9)$$

where $C = C(\text{diam}(\Omega_D), \|\alpha\|_{\mathbf{W}^{1,\infty}(\Omega_D)})$.

Proof. From the error estimates (4.3) and also the analytical estimates (1.16) and (1.18) from discussion of Section 1.1,

$$\begin{aligned} \|f_h - \mathcal{M}\|_{L^2(\Omega_D)} &\leq \|f_h - f\|_{L^2(\Omega_D)} + \|f - \mathcal{M}\|_{L^2(\Omega_D)} \\ &\leq \|f_h - f\|_{L^2(\Omega_D)} + \|f - \mathcal{M}\|_{L^2(\Omega_D)} \\ &\leq \|f_h - f\|_{L^2(\Omega_D)} + \|f - \mathcal{M}\|_{B^2(\mathbb{R}^d)} \\ &\leq C \sqrt{t} e^{Cht} h^{k+\frac{1}{2}} \|f\|_{L^\infty([0,t], H^{k+1}(\Omega_D))} + 3e^{-\lambda t} \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^d)}, \end{aligned}$$

where f^0 is the initial condition. \square

As the above theorem indicates, for a high order discretization with not a large terminal time, we will be able to observe exponential decay of f_h towards equilibrium, which comes in as the term $3e^{-\lambda t} \|f_0 - \mathcal{M}\|_{B^2(\mathbb{R}^d)}$.

4.3. The special case of V_h^0 . In this subsection, we consider the special case of piecewise constant discretization. We show the L^1 stability and positivity-preserving property of the numerical solution as an analog of the exact solution. The following two results are an adaptation of the Crandall-Tartar lemma [25] to low order DG schemes, which states that any mass preserving, contracting linear first order operator is stable and monotone preserving.

THEOREM 13. (L^1 Stability) *Consider the semi-discrete discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9), in the case of piecewise constant basis functions,*

$$\|f_h(t)\|_{L^1(\Omega_{\mathcal{D}})} \leq \|f_h(0)\|_{L^1(\Omega_{\mathcal{D}})}. \quad (4.10)$$

Proof. Fix $t > 0$. Take $w_h = \text{sgn}(\frac{f_h}{M}) = \text{sgn}(f_h)$. This is a valid choice because $w_h \in V_h^0$, and $\text{sgn}(\cdot)$ is also a monotone function. Then we obtain

$$(\partial_t f_h, \text{sgn}(f_h))_{\mathcal{T}_h} + \mathcal{A}(f_h, \text{sgn}(\frac{f_h}{M})) = 0,$$

where

$$\begin{aligned} \mathcal{A}(f_h, \text{sgn}(\frac{f_h}{M})) &= \langle f_h^-, \text{sgn}(f_h) \alpha \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), \text{sgn}(\frac{f_h}{M}))_{\mathcal{T}_h} \\ &= -\frac{1}{2} \langle f_h^-, [\text{sgn}(f_h)] |\alpha \cdot \mathbf{n}| \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), \text{sgn}(\frac{f_h}{M}))_{\mathcal{T}_h} \\ &= \langle |f_h^-|, 1_{\{[\text{sgn}(f_h)] \neq 0\}} |\alpha \cdot \mathbf{n}| \rangle_{\partial \mathcal{T}_h} - (Q_\sigma(f_h), \text{sgn}(\frac{f_h}{M}))_{\mathcal{T}_h} \geq 0 \end{aligned}$$

Integrating over time, (4.10) follows because $(\partial_t f_h, \text{sgn}(f_h))_{\mathcal{T}_h} = \frac{d}{dt} \|f_h(t)\|_{L^1(\Omega_{\mathcal{D}})}$.

\square

THEOREM 14. (Semi-discrete positivity) *Consider the semi-discrete piecewise constant discontinuous Galerkin solution f_h in (3.2) to the linear Boltzmann equation (1.9), provided $f_h(t=0) \geq 0$, the solution remains positive on $\Omega_{\mathcal{D}}$:*

$$f_h(t, x, v) \geq 0, \quad \text{for } t \in [0, T] \text{ for } x, v \in \Omega_{\mathcal{D}}.$$

Proof. The proof of this theorem is rather similar to the previous one and its continuum counterpart. We take the test function $w_h = \frac{1}{2}(\text{sgn}(\frac{f_h}{M}) - 1) = \frac{1}{2}(\text{sgn}(f_h) - 1)$. This is a valid choice because $w_h \in V_h^0$, and $\frac{1}{2}(\text{sgn}(\cdot) - 1)$ is also a monotone function. Then

$$(\partial_t f_h, \frac{1}{2}(\text{sgn}(f_h) - 1))_{\mathcal{T}_h} + \mathcal{A}(f_h, \frac{1}{2}(\text{sgn}(f_h) - 1)) = 0,$$

where

$$\begin{aligned}
\mathcal{A}(f_h, \frac{1}{2}(\text{sgn}(f_h) - 1)) &= \langle f_h^-, \frac{1}{2}(\text{sgn}(f_h) - 1) \alpha \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} - (Q_\sigma(f_h), \frac{1}{2}(\text{sgn}(\frac{f_h}{M}) - 1))_{\mathcal{T}_h} \\
&= -\frac{1}{2} \langle f_h^-, [\frac{1}{2}(\text{sgn}(f_h) - 1)] |\alpha \cdot \mathbf{n}| \rangle_{\partial\mathcal{T}_h} - (Q_\sigma(f_h), \frac{1}{2}(\text{sgn}(\frac{f_h}{M}) - 1))_{\mathcal{T}_h} \\
&= -\frac{1}{4} \langle f_h^-, [\text{sgn}(f_h)] |\alpha \cdot \mathbf{n}| \rangle_{\partial\mathcal{T}_h} - (Q_\sigma(f_h), \frac{1}{2}(\text{sgn}(\frac{f_h}{M}) - 1))_{\mathcal{T}_h} \\
&= \frac{1}{2} \langle |f_h^-|, 1_{\{|\text{sgn}(f_h)| \neq 0\}} |\alpha \cdot \mathbf{n}| \rangle_{\partial\mathcal{T}_h} - (Q_\sigma(f_h), \frac{1}{2}(\text{sgn}(\frac{f_h}{M}) - 1))_{\mathcal{T}_h} \geq 0.
\end{aligned}$$

Let

$$\beta(w) = \begin{cases} 0 & \text{if } w \geq 0 \\ -w & \text{if } w < 0. \end{cases}$$

Then,

$$\frac{d}{dt} \int_{\Omega_{\mathcal{D}}} \beta(f_h) dv dx = (\partial_t f_h, \frac{1}{2}(\text{sgn}(f_h) - 1))_{\mathcal{T}_h} \leq 0.$$

Since the initial condition is assumed to be non-negative, $\beta(f_h(0)) = \max\{0, -f_h(0, x, v)\} = 0$, so that

$$0 \leq \int_{\Omega_{\mathcal{D}}} \beta(f_h(t)) \leq \int_{\Omega_{\mathcal{D}}} \beta(f_h(0)) = 0, \tag{4.11}$$

which implies the non-negativity of $f_h(t)$. \square

5. Analysis of the positivity-preserving DG scheme. In this section, we analyze the positivity-preserving DG scheme proposed in Section 3.3. To simplify the discussion, we will only consider the equation with 1D in x and v space. Similar conclusions hold true for higher dimensions. We present the results in the setting of rectangular meshes. Discussions about non-cartesian meshes are provided after Theorem 15. The discussion below closely follows those in [52], but with collision term treated separately.

Suppose the partition of the computational domain is as follows: $\mathcal{D}_x = \cup_{i=1}^{N_x} I_i$, $\mathcal{D}_v = \cup_{j=1}^{N_v} J_j$, where $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$, $J_j = [v_{j-\frac{1}{2}}, v_{j+\frac{1}{2}}]$. We use Δx_i to denote the length of I_i and Δv_j to denote the length of J_j . Denote an elementary cell K_{ij} as $I_i \times J_j$. We denote by \bar{f}_{ij}^n the cell average of f_h on K_{ij} at time t^n . We abuse the notation and define $f_{i-\frac{1}{2},j}^+(v)$, $f_{i+\frac{1}{2},j}^-(v)$, $f_{i,j-\frac{1}{2}}^+(x)$, $f_{i,j+\frac{1}{2}}^-(x)$ as the traces of f_h on K_{ij} on the four edges at time t^n respectively. Here the $+$ and $-$ is no longer associated with the wind direction α , but the direction of the growth of x, v axis. For example, $f_{i-\frac{1}{2},j}^+(v) = f_h(x_{i-\frac{1}{2}}, v)$ calculated from cell (i, j) , and $f_{i-\frac{1}{2},j}^-(v) = f_h(x_{i-\frac{1}{2}}, v)$ calculated from cell $(i-1, j)$. We now consider the scheme with forward Euler time discretization. TVD-RK methods are convex combinations of Euler forward and are direct generalization of the result below.

We apply the test function $w_h = 1$ on $K_{i,j}$ and $w_h = 0$ elsewhere. One step of Euler forward time discretization of (3.2) will give:

$$\begin{aligned}
\bar{f}_{ij}^{n+1} &= \bar{f}_{ij}^n - \frac{\Delta t}{\Delta x_i \Delta v_j} \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} h_1(f_{i+\frac{1}{2},j}^-(v), f_{i+\frac{1}{2},j}^+(v), v) - h_1(f_{i-\frac{1}{2},j}^-(v), f_{i-\frac{1}{2},j}^+(v), v) dv \\
&\quad - \frac{\Delta t}{\Delta x_i \Delta v_j} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} h_2(f_{i,j+\frac{1}{2}}^-(x), f_{i,j+\frac{1}{2}}^+(x), x) - h_2(f_{i,j-\frac{1}{2}}^-(x), f_{i,j-\frac{1}{2}}^+(x), x) dx \\
&\quad + \frac{\Delta t}{\Delta x_i \Delta v_j} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x, v) dv dx. \tag{5.1}
\end{aligned}$$

In the above formula, $h_1(\cdot, \cdot)$ and $h_2(\cdot, \cdot)$ denote the upwind flux that we have chosen. In particular,

$$\begin{aligned}
h_1(a, b, v) &= v (a \mathbf{1}_{\{v \geq 0\}} + b \mathbf{1}_{\{v < 0\}}), \\
h_2(a, b, x) &= -\frac{e}{m} E(t, x) (a \mathbf{1}_{\{E(t, x) < 0\}} + b \mathbf{1}_{\{E(t, x) \geq 0\}}).
\end{aligned}$$

The integrals in (5.1) can be approximated by quadratures to sufficient accuracy, see Subsection 3.1 for related discussions. Suppose now we use a Gauss quadrature rule with L points, which is exact for single variable polynomials of degree $2L - 1$. From [18], $L \geq k + 1$ will be accurate enough in the sense that the L^∞ norm of error induced by quadrature will be smaller than Ch^{k+1} .

The set of Gauss quadrature points on I_i is defined as $S_i^x = \{x_i^\beta : \beta = 1, \dots, L\}$ and the Gauss quadrature points on J_j is $S_j^v = \{v_j^\beta : \beta = 1, \dots, L\}$. Let w_β be the corresponding Gauss weight on the interval $[-\frac{1}{2}, \frac{1}{2}]$, such that $\sum_\beta w_\beta = 1$. We use subscript β to denote the function value at Gauss quadrature points, for instance, $f_{i+\frac{1}{2},\beta}^- = f_{i+\frac{1}{2},j}^-(v_j^\beta)$, etc. Then (5.1) becomes

$$\begin{aligned}
\bar{f}_{ij}^{n+1} &= \bar{f}_{ij}^n - \frac{\Delta t}{\Delta x_i \Delta v_j} \sum_{\beta=1}^L [h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^\beta) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^\beta)] w_\beta \Delta v_j \\
&\quad - \frac{\Delta t}{\Delta x_i \Delta v_j} \sum_{\beta=1}^L [h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^\beta) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^\beta)] w_\beta \Delta x_i \\
&\quad + \frac{\Delta t}{\Delta x_i \Delta v_j} \sum_{\beta=1}^L \left\{ \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv \right\} w_\beta \Delta x_i. \tag{5.2}
\end{aligned}$$

Denote $\lambda_1 = \frac{\Delta t}{\Delta x_i}$, $\lambda_2 = \frac{\Delta t}{\Delta v_j}$, the above equality becomes,

$$\begin{aligned}
\bar{f}_{ij}^{n+1} &= \bar{f}_{ij}^n - \lambda_1 \sum_{\beta=1}^L w_\beta [h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^\beta) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^\beta)] \\
&\quad - \lambda_2 \sum_{\beta=1}^L w_\beta [h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^\beta) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^\beta)] \\
&\quad + \lambda_2 \sum_{\beta=1}^L w_\beta \left\{ \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv \right\}. \tag{5.3}
\end{aligned}$$

Similar to [52], we write the cell average as

$$\bar{f}_{ij}^n = \sum_{\beta=1}^L w_{\beta} \bar{f}_{\beta j}^n = \sum_{\beta=1}^L w_{\beta} \bar{f}_{i\beta}^n,$$

where $\bar{f}_{\beta j}^n = \frac{1}{\Delta v_j} \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} f_h(x_i^{\beta}, v) dv$, $\bar{f}_{i\beta}^n = \frac{1}{\Delta x_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f_h(x, v_j^{\beta}) dx$. Define $a_1 = \max_{\mathcal{D}_v^D} |v|$, $a_2 = \max_{\mathcal{D}_x^D} |\frac{e}{m} E(t^n, x)|$. Then (5.3) becomes

$$\begin{aligned} \bar{f}_{ij}^{n+1} &= \frac{a_1 \lambda_1}{a_1 \lambda_1 + a_2 \lambda_2} \bar{f}_{ij}^n - \lambda_1 \sum_{\beta=1}^L w_{\beta} [h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^{\beta}) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^{\beta})] \\ &\quad + \frac{a_2 \lambda_2}{a_1 \lambda_1 + a_2 \lambda_2} \bar{f}_{ij}^n - \lambda_2 \sum_{\beta=1}^L w_{\beta} [h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^{\beta}) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^{\beta})] \\ &\quad + \lambda_2 \sum_{\beta=1}^L w_{\beta} \left\{ \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_{\sigma}(f_h)(x_i^{\beta}, v) dv \right\} \\ &= \frac{a_1 \lambda_1}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} \bar{f}_{i\beta}^n - \lambda_1 \sum_{\beta=1}^L w_{\beta} [h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^{\beta}) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^{\beta})] \\ &\quad + \frac{a_2 \lambda_2}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} \bar{f}_{\beta j}^n - \lambda_2 \sum_{\beta=1}^L w_{\beta} [h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^{\beta}) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^{\beta})] \\ &\quad + \lambda_2 \sum_{\beta=1}^L w_{\beta} \left\{ \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_{\sigma}(f_h)(x_i^{\beta}, v) dv \right\} \\ &= \frac{a_1 \lambda_1}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} [\bar{f}_{i\beta}^n - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_1} (h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^{\beta}) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^{\beta}))] \\ &\quad + \frac{a_2 \lambda_2}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} [\bar{f}_{\beta j}^n - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2} (h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^{\beta}) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^{\beta})) \\ &\quad - \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_{\sigma}(f_h)(x_i^{\beta}, v) dv]. \end{aligned}$$

The above equality implies

$$\bar{f}_{ij}^{n+1} = \frac{a_1 \lambda_1}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} H_x^{i,\beta} + \frac{a_2 \lambda_2}{a_1 \lambda_1 + a_2 \lambda_2} \sum_{\beta=1}^L w_{\beta} H_v^{\beta,j},$$

where $H_x^{i,\beta}$, $H_v^{\beta,j}$ denote one dimensional schemes as

$$\begin{aligned} H_x^{i,\beta} &= \bar{f}_{i\beta}^n - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_1} [h_1(f_{i+\frac{1}{2},\beta}^-, f_{i+\frac{1}{2},\beta}^+, v_j^{\beta}) - h_1(f_{i-\frac{1}{2},\beta}^-, f_{i-\frac{1}{2},\beta}^+, v_j^{\beta})], \\ H_v^{\beta,j} &= \bar{f}_{\beta j}^n - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2} [h_2(f_{\beta,j+\frac{1}{2}}^-, f_{\beta,j+\frac{1}{2}}^+, x_i^{\beta}) - h_2(f_{\beta,j-\frac{1}{2}}^-, f_{\beta,j-\frac{1}{2}}^+, x_i^{\beta}) \\ &\quad - \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_{\sigma}(f_h)(x_i^{\beta}, v) dv]. \end{aligned}$$

To have the positivity of cell average at time t^{n+1} , i.e., to ensure $\bar{f}_{ij}^{n+1} \geq 0$, it suffices to have $H_x^{i,\beta} \geq 0$ and $H_v^{\beta,j} \geq 0$. We introduce the Gauss-Lobatto points at this point. We use hats and sub- and super-script α to denote those points, namely, $\hat{S}_i^x = \{\hat{x}_i^\alpha : \alpha = 1, \dots, N\}$ and $\hat{S}_j^v = \{\hat{v}_j^\alpha : \alpha = 1, \dots, N\}$. \hat{w}_α are the Gauss-Lobatto weights on $[-\frac{1}{2}, \frac{1}{2}]$ such that $\sum_\alpha \hat{w}_\alpha = 1$. The N point Gauss-Lobatto quadrature rule will be exact for integrals of polynomials up to degree $(2N-3)$. The number N should be chosen according accuracy of the scheme. For example, if $2N-3 \geq k$, then the equality below is exact.

$$\bar{f}_{i\beta}^n = \sum_{\alpha=1}^N \hat{w}_\alpha f_h(\hat{x}_i^\alpha, v_j^\beta).$$

Define $T_{i,j}^{\alpha,\beta} = f_h(\hat{x}_i^\alpha, v_j^\beta)$ for $\alpha = 1, \dots, N, \beta = 1, \dots, L$ and $T_{i,j}^{0,\beta} = f_h(x_{i-\frac{1}{2}}^-, v_j^\beta) = T_{i-1,j}^{N,\beta}$, $T_{i,j}^{N+1,\beta} = f_h(x_{i+\frac{1}{2}}^+, v_j^\beta) = T_{i+1,j}^{1,\beta}$. Then,

$$\begin{aligned} H_x^{i,\beta} &= \bar{f}_{i\beta}^n - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1} [h_1(T_{i,j}^{N,\beta}, T_{i,j}^{N+1,\beta}, v_j^\beta) - h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)] \\ &= \sum_{\alpha=1}^N \hat{w}_\alpha T_{i,j}^{\alpha,\beta} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1} [h_1(T_{i,j}^{N,\beta}, T_{i,j}^{N+1,\beta}, v_j^\beta) - h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta) \\ &\quad + h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta) - h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)] \\ &= \sum_{\alpha=2}^{N-1} \hat{w}_\alpha T_{i,j}^{\alpha,\beta} + \hat{w}_1 \left[T_{i,j}^{1,\beta} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} [h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta) - h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)] \right] \\ &\quad + \hat{w}_N \left[T_{i,j}^{N,\beta} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_N} [h_1(T_{i,j}^{N,\beta}, T_{i,j}^{N+1,\beta}, v_j^\beta) - h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta)] \right] \end{aligned}$$

We will show that if $T_{i,j}^{\alpha,\beta} \geq 0$ for all $i \in [1, N_x], j \in [1, N_v], \alpha \in [1, N], \beta \in [1, L]$, and $a_1\lambda_1 + a_2\lambda_2 \leq \min(\hat{w}_1, \hat{w}_N)$, then $H_x^{i,\beta} \geq 0$. Define $Q_{i,j} = T_{i,j}^{1,\beta} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} [h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta) - h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)]$, then this is a monotonically increasing function for all of its variables $T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}$, which can be verified by computing its partial derivatives.

$$\frac{\partial Q_{i,j}}{\partial T_{i,j}^{0,\beta}} = \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} \frac{\partial h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)}{\partial T_{i,j}^{0,\beta}} = \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} v_j^\beta 1_{\{v_j^\beta \geq 0\}} \geq 0,$$

$$\frac{\partial Q_{i,j}}{\partial T_{i,j}^{N,\beta}} = -\frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} \frac{\partial h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta)}{\partial T_{i,j}^{N,\beta}} = -\frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} v_j^\beta 1_{\{v_j^\beta < 0\}} \geq 0,$$

$$\begin{aligned} \frac{\partial Q_{i,j}}{\partial T_{i,j}^{1,\beta}} &= 1 - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} \left(\frac{\partial h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta)}{\partial T_{i,j}^{1,\beta}} - \frac{\partial h_1(T_{i,j}^{0,\beta}, T_{i,j}^{1,\beta}, v_j^\beta)}{\partial T_{i,j}^{1,\beta}} \right) \\ &= 1 - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} \left(v_j^\beta 1_{\{v_j^\beta \geq 0\}} - v_j^\beta 1_{\{v_j^\beta < 0\}} \right) \\ &= 1 - \frac{a_1\lambda_1 + a_2\lambda_2}{a_1\hat{w}_1} |v_j^\beta| \geq 1 - \frac{a_1\lambda_1 + a_2\lambda_2}{\hat{w}_1} \geq 0. \end{aligned}$$

The last line of inequality is because $a_1 = \max_{\mathcal{D}_v} |v| \geq |v_j^\beta|$. Since $T_{i,j}^{0,\beta} \geq 0$, $T_{i,j}^{1,\beta} \geq 0$, $T_{i,j}^{N,\beta} \geq 0$,

$$Q_{i,j} \geq 0 - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_1 \hat{w}_\alpha} [h_1(0, 0, v_j^\beta) - h_1(0, 0, v_j^\beta)] = 0.$$

Similarly, we can deduce, $T_{i,j}^{N,\beta} - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_1 \hat{w}_N} [h_1(T_{i,j}^{N,\beta}, T_{i,j}^{N+1,\beta}, v_j^\beta) - h_1(T_{i,j}^{1,\beta}, T_{i,j}^{N,\beta}, v_j^\beta)] \geq 0$. Hence, $H_x^{i,\beta} \geq 0$.

In summary, we have proved that if $T_{i,j}^{\alpha,\beta} = f_h(\hat{x}_i^\alpha, v_j^\beta) \geq 0$ for all $i \in [1, N_x]$, $j \in [1, N_v]$, $\alpha \in [1, N]$, $\beta \in [1, L]$, and the CFL condition $a_1 \lambda_1 + a_2 \lambda_2 \leq \min(\hat{w}_1, \hat{w}_N)$ is satisfied, then $H_x^{i,\beta} \geq 0$. In the derivation, we only require $2N - 3 \geq k$. One can show easily that for the nontrivial case $k \geq 1$, $N = k + 1$ satisfies this condition. From this point on, we assume $N = k + 1$.

Similar argument will follow for $H_v^{\beta,j}$, with the difference coming from the collisional integral $\int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv$. In fact,

$$\int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv = \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{\mathcal{D}_v} (\sigma(x_i^\beta, v, v') f_h(x_i^\beta, v') - \sigma(x_i^\beta, v', v) f_h(x_i^\beta, v)) dv' dv$$

Since $f_h(x_i^\beta, v) \in P^k(v)$ on cell J_j , it can be expanded as a linear combination of the basis functions

$$f_h(x_i^\beta, v) = \sum_{\alpha=1}^N f_h(x_i^\beta, \hat{v}_j^\alpha) L_\alpha\left(\frac{v - v_j}{\Delta v_j}\right) \quad \text{on } J_j,$$

where $L_\alpha(\cdot)$ are basis functions for Gauss-Lobatto points s_α on the interval $[-\frac{1}{2}, \frac{1}{2}]$ such that $L_\alpha(s_\gamma) = \delta_{\alpha\gamma}$. Hence,

$$\begin{aligned} & \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv = \\ & \sum_{m=1}^{N_v} \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} (\sigma(x_i^\beta, v, v') f_h(x_i^\beta, v') - \sigma(x_i^\beta, v', v) f_h(x_i^\beta, v)) dv' dv \\ & = \sum_{m=1}^{N_v} \sum_{\alpha=1}^N \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} (\sigma(x_i^\beta, v, v') f_h(x_i^\beta, \hat{v}_m^\alpha) L_\alpha\left(\frac{v' - v_m}{\Delta v_m}\right) \\ & \quad - \sigma(x_i^\beta, v', v) f_h(x_i^\beta, \hat{v}_j^\alpha) L_\alpha\left(\frac{v - v_j}{\Delta v_j}\right)) dv' dv \\ & = \sum_{\alpha=1}^N \sum_{m=1}^{N_v} [f_h(x_i^\beta, \hat{v}_m^\alpha) \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} \sigma(x_i^\beta, v, v') L_\alpha\left(\frac{v' - v_m}{\Delta v_m}\right) dv' dv \\ & \quad - f_h(x_i^\beta, \hat{v}_j^\alpha) \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} \sigma(x_i^\beta, v', v) L_\alpha\left(\frac{v - v_j}{\Delta v_j}\right) dv' dv] \end{aligned}$$

If we define $A_{j,m}^{\alpha,\beta} = \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} \sigma(x_i^\beta, v, v') L_\alpha\left(\frac{v' - v_m}{\Delta v_m}\right) dv' dv$, then it can be verified that

$$\int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} Q_\sigma(f_h)(x_i^\beta, v) dv = \sum_{\alpha=1}^N \sum_{m=1}^{N_v} \{ f_h(x_i^\beta, \hat{v}_m^\alpha) A_{j,m}^{\alpha,\beta} - f_h(x_i^\beta, \hat{v}_j^\alpha) A_{m,j}^{\alpha,\beta} \}.$$

Here, $A_{j,m}^{\alpha,\beta}$ satisfies the following important properties.

(1) If we denote the bound for the basis functions as $L_\alpha(\cdot) \leq C_k$ on $[-\frac{1}{2}, \frac{1}{2}]$, where C_k is a constant that only depends on the polynomial order k , then

$$\begin{aligned} \sum_{m=1}^{N_v} A_{j,m}^{\alpha,\beta} &= \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{\mathcal{D}_v} \sigma(x_i^\beta, v, v') L_\alpha\left(\frac{v' - v_m}{\Delta v_m}\right) dv' dv \\ &\leq C_k \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \int_{\mathcal{D}_v} \sigma(x_i^\beta, v, v') dv' dv \leq C_k K \Delta v_j, \end{aligned}$$

where we have used Property 1 of the collision integral.

(2) For piecewise linear case, i.e. when $k = 1$, since $L_1(\cdot)$, $L_2(\cdot)$, $\sigma(\cdot, \cdot, \cdot)$ are always positive, we have $A_{j,m}^{\alpha,\beta} \geq 0$ for $\forall i, j, \alpha, \beta$. For the general case of $k > 1$, the positivity of $A_{j,m}^{\alpha,\beta}$ cannot be guaranteed. However, if we use the same Gauss-Lobatto quadrature to evaluate the integration in the variable v' , we have

$$A_{j,m}^{\alpha,\beta} \approx \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} \sigma(x_i^\beta, v, \hat{v}_m^\alpha) \hat{w}_\alpha dv \geq 0.$$

This quadrature corresponds to the requirement that when evaluating the collisional integral term in the DG scheme, in the v and v' variable we use a $N = k + 1$ point Gauss-Lobatto quadrature. This is a reasonable assumption, see for instance, the discussion at the end of section 3.1. One can show that by doing this, an error of at most order $h^{k+\frac{1}{2}}$ will be introduced to the solution, see [18]. On the other hand, in the very special case of relaxation models, which will be presented in Section 6, this restriction can be lifted. Because $\sigma(x, v, v') = kM(v)$, and

$$A_{j,m}^{\alpha,\beta} = k \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} M(v) dv \cdot \int_{v_{m-\frac{1}{2}}}^{v_{m+\frac{1}{2}}} L_\alpha\left(\frac{v' - v_m}{\Delta v_m}\right) dv' = k \int_{v_{j-\frac{1}{2}}}^{v_{j+\frac{1}{2}}} M(v) dv \geq 0.$$

From here on, we will assume $A_{j,m}^{\alpha,\beta} \geq 0$ for $\forall i, j, \alpha, \beta$.

Now we return to $H_v^{\beta,j}$, since

$$\bar{f}_{\beta j}^n = \sum_{\alpha=1}^N \hat{w}_\alpha f_h(x_i^\beta, \hat{v}_j^\alpha).$$

Define $S_{i,j}^{\beta,\alpha} = f_h(x_i^\beta, \hat{v}_j^\alpha)$ for $\alpha = 1, \dots, N, \beta = 1, \dots, L$ and $S_{i,j}^{\beta,0} = f_h(x_i^\beta, v_{j-\frac{1}{2}}^-) = S_{i,j-1}^{\beta,N}$, $S_{i,j}^{\beta,N+1} = f_h(x_i^\beta, v_{j+\frac{1}{2}}^+) = S_{i,j+1}^{\beta,1}$. Then,

$$\begin{aligned} H_v^{\beta,j} &= \bar{f}_{\beta j}^n - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2} [h_2(S_{i,j}^{\beta,N}, S_{i,j}^{\beta,N+1}, x_i^\beta) - h_2(S_{i,j}^{\beta,0}, S_{i,j}^{\beta,1}, x_i^\beta) \\ &\quad - \sum_{\alpha=1}^N \sum_{m=1}^{N_v} (S_{i,m}^{\beta,\alpha} A_{j,m}^{\alpha,\beta} - S_{i,j}^{\beta,\alpha} A_{m,j}^{\alpha,\beta})] \\ &= \sum_{\alpha=2}^{N-1} \hat{w}_\alpha \mathcal{H}_{i,j}^{\beta,\alpha} + \hat{w}_1 \mathcal{F}_{i,j} + \hat{w}_N \mathcal{G}_{i,j}, \end{aligned}$$

where

$$\begin{aligned}\mathcal{F}_{i,j} &= S_{i,j}^{\beta,1} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} \{h_2(S_{i,j}^{\beta,1}, S_{i,j}^{\beta,N}, x_i^\beta) \\ &\quad - h_2(S_{i,j}^{\beta,0}, S_{i,j}^{\beta,1}, x_i^\beta) - \sum_{m=1}^{N_v} (S_{i,m}^{\beta,1} A_{j,m}^{1,\beta} - S_{i,j}^{\beta,1} A_{m,j}^{1,\beta})\},\end{aligned}$$

$$\begin{aligned}\mathcal{G}_{i,j} &= S_{i,j}^{\beta,N} - \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_N} \{h_2(S_{i,j}^{\beta,N}, S_{i,j}^{\beta,N+1}, x_i^\beta) \\ &\quad - h_2(S_{i,j}^{\beta,1}, S_{i,j}^{\beta,N}, x_i^\beta) - \sum_{m=1}^{N_v} (S_{i,m}^{\beta,N} A_{j,m}^{N,\beta} - S_{i,j}^{\beta,N} A_{m,j}^{N,\beta})\},\end{aligned}$$

and

$$\mathcal{H}_{i,j}^{\beta,\alpha} = S_{i,j}^{\beta,\alpha} + \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_\alpha} \sum_{m=1}^{N_v} (S_{i,m}^{\beta,\alpha} A_{j,m}^{\alpha,\beta} - S_{i,j}^{\beta,\alpha} A_{m,j}^{\alpha,\beta}).$$

Similar to the argument before, to have $H_v^{\beta,j} \geq 0$, we only need $\mathcal{F}_{i,j} \geq 0$, $\mathcal{G}_{i,j} \geq 0$, $\mathcal{H}_{i,j}^{\beta,\alpha} \geq 0$. We will prove that this will hold if $S_{i,j}^{\beta,\alpha} = f_h(x_i^\beta, \hat{v}_j^\alpha) \geq 0$ for all $i \in [1, N_x]$, $j \in [1, N_v]$, $\alpha \in [1, N]$, $\beta \in [1, L]$, and a proper CFL condition is true. For simplicity, we will only prove for $\mathcal{F}_{i,j} \geq 0$. The other two terms follow under similar argument. In particular, we will show $\mathcal{F}_{i,j}$ is an increasing function of all its variables.

$$\begin{aligned}\frac{\partial \mathcal{F}_{i,j}}{\partial S_{i,j}^{\beta,0}} &= \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} \frac{\partial h_2(S_{i,j}^{\beta,0}, S_{i,j}^{\beta,1}, x_i^\beta)}{\partial S_{i,j}^{\beta,0}} \\ &= \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} \left(-\frac{e}{m} E(t, x_i^\beta)\right) 1_{\{E(t, x_i^\beta) \leq 0\}} \geq 0,\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{F}_{i,j}}{\partial S_{i,j}^{\beta,N}} &= -\frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} \frac{\partial h_2(S_{i,j}^{\beta,1}, S_{i,j}^{\beta,N}, x_i^\beta)}{\partial S_{i,j}^{\beta,N}} \\ &= \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} \frac{e}{m} E(t, x_i^\beta) 1_{\{E(t, x_i^\beta) > 0\}} \geq 0,\end{aligned}$$

$$\frac{\partial \mathcal{F}_{i,j}}{\partial S_{i,m}^{\beta,1}} = \frac{a_1\lambda_1 + a_2\lambda_2}{a_2\hat{w}_1} A_{j,m}^{1,\beta} \geq 0 \quad \text{if } m \neq j,$$

and

$$\begin{aligned}
\frac{\partial \mathcal{F}_{i,j}}{\partial S_{i,j}^{\beta,1}} &= 1 - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2 \hat{w}_1} \times \\
&\quad \left(\frac{\partial h_2(S_{i,j}^{\beta,1}, S_{i,j}^{\beta,N}, x_i^\beta)}{\partial S_{i,j}^{\beta,1}} - \frac{\partial h_2(S_{i,j}^{\beta,0}, S_{i,j}^{\beta,1}, x_i^\beta)}{\partial S_{i,j}^{\beta,1}} - A_{j,j}^{1,\beta} + \sum_{m=1}^{N_v} A_{m,j}^{1,\beta} \right) \\
&= 1 - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2 \hat{w}_1} \left[-\frac{e}{m} E(t, x_i^\beta) 1_{\{E(t, x_i^\beta) \leq 0\}} + \frac{e}{m} E(t, x_i^\beta) 1_{\{E(t, x_i^\beta) > 0\}} \right. \\
&\quad \left. - A_{j,j}^{1,\beta} + \sum_{m=1}^{N_v} A_{m,j}^{1,\beta} \right] \\
&= 1 - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2 \hat{w}_1} \left[\left| \frac{e}{m} E(t, x_i^\beta) \right| - A_{j,j}^{1,\beta} + \sum_{m=1}^{N_v} A_{m,j}^{1,\beta} \right] \\
&\geq 1 - \frac{a_1 \lambda_1 + a_2 \lambda_2}{a_2 \hat{w}_1} [a_2 + C_k K \Delta v_j] \geq 0.
\end{aligned}$$

For the last inequality to hold, we need a slightly more restrictive CFL condition, i.e.

$$a_1 \lambda_1 + a_2 \lambda_2 \leq \frac{\hat{w}_1}{1 + \frac{C_k K \max_j \Delta v_j}{a_2}}.$$

The other CFL restrictions derived from $\mathcal{G}_{i,j} \geq 0$, $\mathcal{H}_{i,j}^{\beta,\alpha} \geq 0$ will be

$$a_1 \lambda_1 + a_2 \lambda_2 \leq \frac{\hat{w}_N}{1 + \frac{C_k K \max_j \Delta v_j}{a_2}},$$

and

$$a_1 \lambda_1 + a_2 \lambda_2 \leq \frac{a_2 \min_{\alpha=2,\dots,N-1} \hat{w}_\alpha}{C_k K \max_j \Delta v_j}.$$

When $\max_j \Delta v_j \rightarrow 0$, the CFL condition approaches to $a_1 \lambda_1 + a_2 \lambda_2 \leq \min(\hat{w}_1, \hat{w}_N)$. For finite mesh size, we can assume, for example, when $\max_j \Delta v_j \leq \frac{a_2}{s C_k K}$, where s is a constant, it is enough to have $a_1 \lambda_1 + a_2 \lambda_2 \leq \frac{1}{2} \min(\hat{w}_1, \hat{w}_N)$.

Our conclusion is, if $S_{i,j}^{\beta,\alpha} = f_h(x_i^\beta, \hat{v}_j^\alpha) \geq 0$ for all $i \in [1, N_x], j \in [1, N_v], \alpha \in [1, N], \beta \in [1, L]$, and the CFL condition $a_1 \lambda_1 + a_2 \lambda_2 \leq \min(\hat{w}_1, \hat{w}_N)/2$ is satisfied, then $H_v^{\beta,j} \geq 0$. Put those two parts together, we have our main theorem in this section.

THEOREM 15. *Consider the semi-discrete DG scheme of piecewise P^k polynomials with forward Euler time stepping on a rectangular mesh that is refined enough, if at time t^n , we have $f_h(x, v) \geq 0$ on the set $S_{i,j} = (S_i^x \otimes \hat{S}_j^y) \cup (\hat{S}_i^x \otimes S_j^y)$ for all i, j , where S and \hat{S} denote Gauss and Gauss-Lobatto quadrature points with $L \geq k+1$ and $N = k+1$ points respectively. Moreover, if the CFL condition $a_1 \lambda_1 + a_2 \lambda_2 \leq \min(\hat{w}_1, \hat{w}_N)/2$ is satisfied, then we have the cell average at the next time step t^{n+1} will be positive, i.e.*

$$\bar{f}_{ij}^{n+1} \geq 0 \quad \text{for all } i, j.$$

If $k > 1$, then we require that N point Gauss-Lobatto quadrature rules when evaluating the collision term. This restriction can be lifted for the relaxation model.

Remark: It is well known that for conservation laws, forward Euler time stepping coupled with DG schemes of P^k with $k \geq 1$ under the CFL condition described in the above theorem is unconditionally unstable [24]. We only use this theorem as an intermediate step to the applications with TVD-RK time discretizations, which are convex combinations of forward Euler schemes.

The positivity of the f_h on set $S_{i,j}$ will be guaranteed by the limiter we proposed in Section 3.3. We repeat the definition here: before each forward Euler time step in the TVD-RK scheme,

- on each cell K_{ij} , evaluate $T_{ij} = \min_{(x,v) \in S_{i,j}} f_h(x, v)$.
- Compute $\tilde{f}_h(x, v) = \theta(f_h(x, v) - \bar{f}_{ij}) + \bar{f}_{ij}$, where $\theta = \min\{1, |\bar{f}_{ij}|/|T_{ij} - \bar{f}_{ij}|\}$.
- Use \tilde{f}_h instead of f_h to compute the Euler forward step. A proper CFL condition needs to be enforced as mentioned in Theorem 15.

Remark: The positivity-preserving DG scheme in this paper will ensure the cell averages of the numerical solutions f_h is positive, and \tilde{f}_h is positive on the set of $S_{i,j}$. However, it does not guarantee that f_h, \tilde{f}_h will be positive at every point on the computational domain. For practical purposes, this type of positivity is enough. A stronger limiter by using $T_{ij} = \min_{(x,v) \in K_{i,j}} f_h(x, v)$ can guarantee that \tilde{f}_h is positive on all points in $\Omega_{\mathcal{D}}$ [51]. If one choose to use the strong limiter with $T_{ij} = \min_{(x,v) \in K_{i,j}} f_h(x, v)$ at the last time step, one can also recover pointwise positivity. We also want to remark if this type of limiter is enforced in place of the one proposed in Theorem 15, then the restriction of using quadratures points to evaluate the collision term can be completely removed. This is because when f_h is pointwise positive, the gain term in the collision operator will be positive automatically. Hence, in the proof, the requirement of $A_{j,m}^{\alpha,\beta} \geq 0$ is unnecessary. We also remark that pointwise positivity together with mass conservation will imply L^1 stability of the fully-discrete scheme. In spite of the above facts, this limiter is more difficult to implement for high order schemes in high dimensions, because it requires finding the minimum values of a function in each cell.

In [52], it is proven that this limiter will not destroy the accuracy of the scheme. In particular, it only introduces error in the L^∞ norm of order h^{k+1} . There might be some small accuracy loss due to the TVD-RK or SSP-RK time discretizations, but it is not significant and can be overcome by using a SSP multi-step time stepping, see [52] for details. This limiter keeps the cell average of f_h , so it will not destroy mass conservation. The computational cost induced by the limiter is very small and the CFL condition is on the same order of the one without the limiter. Because of the above reasons, to have a physically-relevant solution, the limiters and the positivity-preserving DG scheme in this section should be used to avoid negative *pdfs* especially when the computations of quantities such as entropies are desired. The implementation of this limiter on non-cartesian meshes is more involved, especially in terms of the location of quadrature points and the treatment of collision operator. A recent paper [54] has relevant discussions on this type of issues. We will explore this aspect in our future work.

6. Numerical Results. We present a numerical implementation and results for the one-dimensional in physical and velocity space test problem for the relaxation model (1.5), in a suitable cut-off rectangular domain $\Omega_{\mathcal{D}} = \mathcal{D}_x^1 \times \mathcal{D}_v^1 = [-L, L] \times$

$[-V, V]$. In particular, we consider

$$\begin{aligned}
\frac{\partial f}{\partial t} + \alpha \cdot \nabla f &= \frac{1}{\tau} (M_\theta(v) \rho(t, x) - f), \quad (x, v) \in \Omega_{\mathcal{D}}, t \in \mathbb{R}_+ \\
f(t, -L, v) &= 0, \quad \text{for } x = -L, 0 < v \leq V, \\
f(t, L, v) &= 0, \quad \text{for } x = L, -V \leq v < 0, \\
f(t, x, v) &= 0, \quad \text{for } v = +V, -V, x \in D_x^1 \\
f(0, x, v) &= f_0(x, v),
\end{aligned} \tag{6.1}$$

with constant collision frequency $\tau^{-1} = 1$ and $\theta = 1$. We take a cut off domain $\mathcal{D}_v^1 = [-V, V] = [-5, 5]$, which are sufficiently large to assume homogeneous Dirichlet boundary conditions for f , since $\epsilon = 1 - \int_{-5}^5 \mu_\infty(x, v) \approx e^{-20}$ for $\theta = 1$ and $\mu_\infty(\pm 5) \approx e^{-25}$, so setting $f(t, x, \pm 5) = 0$ will not produce any mass loss larger than ϵ for a quadratically confined potential $V(x)$, in accordance with the long time behavior of the solution of the kinetic equation in the whole (x, v) -space. Also $\mathcal{D}_x^1 = [-L, L] = [-5, 5]$.

6.1. Comparison of the traditional RKDG scheme with the positivity-preserving DG scheme. In this subsection, we will compare the traditional RKDG scheme with the positivity-preserving DG scheme with special emphasis put on accuracy and positivity of the numerical solution. In particular, we consider (6.1) with potential $V(x) = -x^2/2$ and initial condition $f_0 = \frac{1}{s} \sin(x^2/2)^2 \exp(-(x^2 + v^2)/2)$, where s is the normalizing constant, such that $\int_{\mathcal{D}_v} \int_{\mathcal{D}_x} f_0 dx dv = 1$. Because of the lack of the exact solution, in Table 6.1, we present the difference of DG schemes with and without the positivity-preserving limiter (not the real errors compared to the exact solution). The computation is performed under the same number of mesh point in all directions until the same final state for both methods. We observe the order of convergence for the difference of two approach is optimal despite the degeneracy of accuracy for SSP-RK time discretization reported in [52].

TABLE 6.1

The difference of DG schemes with and without the positivity-preserving limiter when using P^2 polynomials and third order RK time discretization until $T = 0.1$. CFL = 0.2

Mesh	L^1 difference	L^1 order	L^∞ difference	L^∞ order
25 X 25	1.98E-5	-	1.04E-3	-
50 X 50	1.14E-6	4.11	1.46E-4	2.82
100 X 100	8.95E-8	3.67	2.00E-5	2.87
200 X 200	3.40E-9	4.72	6.63E-7	4.91

In Figure 6.1, we plot for the DG scheme without the limiter, the evolution of the number of cells with negative cell averages. We have verified that the positivity preserving DG scheme does not produce any cells with negative averages or cell center values at any time step after $t = 0$.

6.2. Comparisons of decay rates to equilibrium, measured by different entropy functionals. In this subsection, we consider the decay rate of an initial state to equilibrium. In particular, we consider

$$\mathcal{H}_{\log}(t) = \int_{\Omega_{\mathcal{D}}} H \log H \mathcal{M}(x, v) dx dv, \tag{6.2}$$

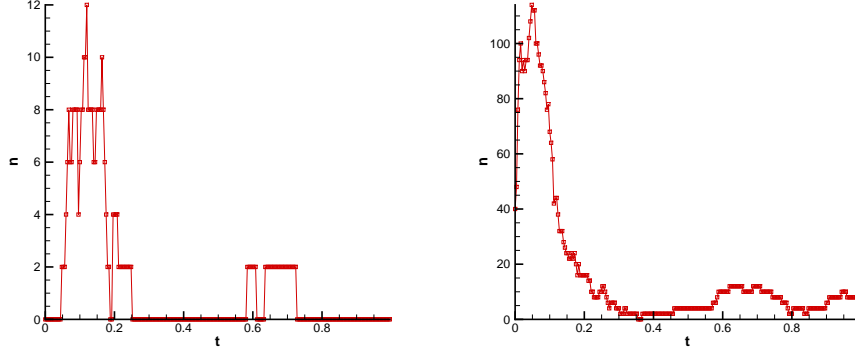


FIG. 6.1. The number of cells (out of 2500 cells) with negative cell averages (left figure) and negative cell center values (right figure) as a function of time for the traditional DG scheme. The computation is performed on a 50×50 mesh with piecewise quadratic polynomials and third order Runge-Kutta time stepping.

that measures (what we refer to as) $f \log f$ -decay in time and the quadratic H functional

$$\mathcal{H}_2(t) = \int_{\Omega_D} H^2 \mathcal{M} dx dv. \quad (6.3)$$

Here \mathcal{M} is the equilibrium distribution, and $H(t, x, v) = f_h(t, x, v)/\mathcal{M}(x, v)$ is the global relative entropy function. We remark that the computation of \mathcal{H}_{\log} requires the positivity of the function $f_h(t, x, v)$. This can be guaranteed by using the positivity-preserving DG schemes coupled with the $(k+1)$ -th order accurate limiter described in [51] and below.

Let $T_K = \min_{(x,v) \in K} f_h(x, v)$ and compute $\tilde{f}_h(x, v) = \theta(f_h(x, v) - \bar{f}_K) + \bar{f}_K$, where $\theta = \min\{1, |\bar{f}_K|/|T_K - \bar{f}_K|\}$. Then $\tilde{f}_h(x, v) \geq 0$ everywhere.

Our numerical examples and simulations are computed under the same initial condition and potential as the previous subsection. In Figure 6.2, we plot the two decay rates. In Figure 6.3, we show the evolution of pdf towards the equilibrium distribution.

7. Conclusions. In this paper, we develop a high-order positivity-preserving DG scheme for linear Vlasov-Boltzmann transport equations (BTE) under the action of quadratically confined electrostatic potentials. Future work includes generalization of the solver to higher dimensions on arbitrary triangulations, and to nonlinear Boltzmann-Poisson system in semiconductor device simulations, especially for electron-hole transport.

Acknowledgments. The authors thank Chi-Wang Shu and Xiangxiong Zhang for discussions about the maximum-principle-satisfying schemes for conservation laws. Yingda Cheng has been partially funded by Focus Research Grant-0757450 and grant NSF DMS-1016001. Irene M. Gamba is supported by grant NSF DMS-0807712 and also DMS-0757450. Support from the Institute of Computational Engineering and Sciences and the University of Texas Austin is gratefully acknowledged.

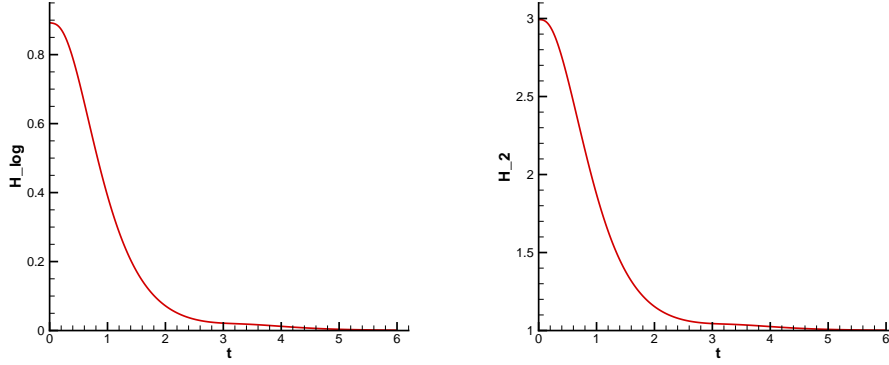


FIG. 6.2. Decay rate. Left: \mathcal{H}_{\log} , right: \mathcal{H}_2 . Positivity-preserving DG scheme computed on a 50×50 mesh with piecewise quadratic polynomials and third order Runge-Kutta time stepping.

REFERENCES

- [1] V. R. A.M. Anile, G. Mascali. Recent developments in hydrodynamical modeling of semiconductors, mathematical problems in semiconductor physics. *Lect. Notes Math.*, 1823:156.
- [2] N. Ben Abdallah and M. L. Tayeb. Diffusion approximation for the one dimensional Boltzmann-Poisson system. *Discrete Contin. Dyn. Syst. Ser. B*, 4(4):1129–1142, 2004.
- [3] G. A. Bird. *Molecular gas dynamics*. Clarendon Press, Oxford, 1994.
- [4] J. E. Broadwell. Study of the rarified shear flow by the discrete velocity method. *J. Fluid Mech.*, 19:401–414, 1964.
- [5] M. J. Caceres, J. A. Carrillo, I. M. Gamba, A. Majorana, and C.-W. Shu. Deterministic kinetic solvers for charged particle transport in semiconductor devices. *Transport Phenomena and Kinetic Theory Applications to Gases, Semiconductors, Photons, and Biological Systems*, pages 151–171, 2007.
- [6] J. A. Carrillo, I. M. Gamba, A. Majorana, and C.-W. Shu. A WENO-solver for 1D non-stationary Boltzmann-Poisson system for semiconductor devices. *J. Comput. Electron.*, 1:365–375, 2002.
- [7] J. A. Carrillo, I. M. Gamba, A. Majorana, and C.-W. Shu. A direct solver for 2D non-stationary Boltzmann-Poisson systems for semiconductor devices: a MESFET simulation by WENO-Boltzmann schemes. *J. Comput. Electron.*, 2:375–380, 2003.
- [8] J. A. Carrillo, I. M. Gamba, A. Majorana, and C.-W. Shu. A WENO-solver for the transients of Boltzmann-Poisson system for semiconductor devices. performance and comparisons with Monte Carlo methods. *J. Comput. Phys.*, 184:498–525, 2003.
- [9] J. A. Carrillo, I. M. Gamba, A. Majorana, and C.-W. Shu. 2D semiconductor device simulations by WENO-Boltzmann schemes: efficiency, boundary conditions and comparison to Monte Carlo methods. *J. Comput. Phys.*, 214:55–80, 2006.
- [10] Y. Cheng, I. M. Gamba, A. Majorana, and C.-W. Shu. Discontinuous Galerkin solver for the semiconductor Boltzmann equation. *SISPAD 07, June 14-17*, pages 257–260, 2007.
- [11] Y. Cheng, I. M. Gamba, A. Majorana, and C.-W. Shu. Discontinuous Galerkin solver for Boltzmann-Poisson transients. *J. Comput. Electron.*, 7:119–123, 2008.
- [12] Y. Cheng, I. M. Gamba, A. Majorana, and C.-W. Shu. A Discontinuous Galerkin solver for Boltzmann-Poisson systems for semiconductor devices. *Comput. Methods Appl. Mech. Eng.*, 198:3130–3150, 2009.
- [13] Y. Cheng, I. M. Gamba, A. Majorana, and C.-W. Shu. A discontinuous Galerkin solver for full-band Boltzmann-Poisson models. *Proceeding of the IWCE13*, pages 211–214, 2009.
- [14] Y. Cheng, I. M. Gamba, A. Majorana, and C.-W. Shu. Performance of discontinuous Galerkin solver for semiconductor Boltzmann Equation. *Proceedings of the IWCE14*, page to appear., 2010.
- [15] P. Ciarlet. *The finite element method for elliptic problems*. North-Holland, Amsterdam, 1975.
- [16] B. Cockburn, B. Dong, and J. Guzmán. Optimal convergence of the original DG method for the transport-reaction equation on special meshes. *SIAM J. Numer. Anal.*, 46:1250–1265,

- 2008.
- [17] B. Cockburn, B. Dong, J. Guzmán, and J. Qian. Optimal convergence of the original dg method in special meshes for variable velocity. 2009. preprint.
 - [18] B. Cockburn, S. Hou, and C.-W. Shu. The Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws IV: the multidimensional case. *Math. Comput.*, 54:545–581, 1990.
 - [19] B. Cockburn, G. Kanschat, I. Perugia, and D. Schötzau. Superconvergence of the local discontinuous Galerkin method for elliptic problems on Cartesian grids. *SIAM J. Numer. Anal.*, 39:264–285, 2001.
 - [20] B. Cockburn, S. Y. Lin, and C.-W. Shu. TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws III: one dimensional systems. *J. Comput. Phys.*, 84:90–113, 1989.
 - [21] B. Cockburn and C.-W. Shu. TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws II: general framework. *Math. Comput.*, 52:411–435, 1989.
 - [22] B. Cockburn and C.-W. Shu. The Runge-Kutta local projection p1-discontinuous Galerkin finite element method for scalar conservation laws. *Math. Model. Num. Anal.*, 25:337–361, 1991.
 - [23] B. Cockburn and C.-W. Shu. The Runge-Kutta discontinuous Galerkin method for conservation laws V: multidimensional systems. *J. Comput. Phys.*, 141:199–224, 1998.
 - [24] B. Cockburn and C.-W. Shu. Runge-Kutta discontinuous Galerkin methods for convection-dominated problems. *J. Sci. Comput.*, 16:173–261, 2001.
 - [25] M. G. Crandall and L. Tartar. Some relations between nonexpansive and order preserving mappings. *Proc. Amer. Math. Soc.*, 78:385–390, 1980.
 - [26] E. Fatemi and F. Odeh. Upwind finite difference solution of Boltzmann equation applied to electron transport in semiconductor devices. *J. Comput. Phys.*, 108:209–217, 1993.
 - [27] I. Gamba and S. H. Tharkabhushaman. Spectral-Lagrangian based methods applied to computation of non-equilibrium statistical states. *J. Comput. Phys.*, 228:2012–2036, 2009.
 - [28] S. Gottlieb and C.-W. Shu. Total variation diminishing Runge-Kutta schemes. *Math. Comput.*, 67:73–85, 1998.
 - [29] S. Gottlieb, C.-W. Shu, and E. Tadmor. Strong stability preserving high order time discretization methods. *SIAM Review*, 43:89–112, 2001.
 - [30] Y. Guo. The Vlasov-Poisson-Boltzmann system near Maxwellians. *Comm. Pure Appl. Math.*, 55(9):1104–1135, 2002.
 - [31] Y. Guo. The Vlasov-Maxwell-Boltzmann system near Maxwellians. *Invent. Math.*, 153(3):593–630, 2003.
 - [32] B. Helffer and F. Nier. Hypocoelliptic estimates and spectral theory for Fokker-Planck operators and Witten laplacians. *Lecture Notes in Mathematics series*, 1862, 2005.
 - [33] F. Hérau. Hypocoercivity and exponential time decay for the linear inhomogeneous relaxation Boltzmann equation. *Asymptot. Anal.*, 46:349–359, 2006.
 - [34] F. Hérau and F. Nier. Isotropic hypocoercivity and trend to equilibrium for the Fokker-Planck equation with high degree potential. *Arch. Ration. Mech. Anal.*, 171:151–218, 2004.
 - [35] C. Johnson and J. Pitkäranta. An analysis of the discontinuous Galerkin method for a scalar hyperbolic equation. *Math. Comp.*, 46:1–26, 1986.
 - [36] P. Lesaint and P.-A. Raviart. On a finite element method for solving the neutron transport equation. In *Mathematical aspects of finite elements in partial differential equations (Proc. Sympos., Math. Res. Center, Univ. Wisconsin, Madison, Wis., 1974)*, pages 89–123. Math. Res. Center, Univ. of Wisconsin-Madison, Academic Press, New York, 1974.
 - [37] A. Majorana and R. Piatella. A finite difference scheme solving the Boltzmann Poisson system for semiconductor devices. *J. Comput. Phys.*, 174:649–668, 2001.
 - [38] C. Mouhot and L. Neumann. Quantitative perturbative study of convergence to equilibrium for collisional kinetic models in the torus. *Nonlinearity*, 4:969–998, 2006.
 - [39] K. Nanbu. Direct simulation scheme derived from the Boltzmann equation. i. monocomponent gases. *J. Phys. Soc. Jpn.*, 52:2042–2049, 1983.
 - [40] L. Pareschi and B. Perthame. A Fourier spectral method for homogeneous Boltzmann equations. *Transport Theory Statist. Phys.*, 25:369C383, 1996.
 - [41] T. E. Peterson. A note on the convergence of the discontinuous Galerkin method for a scalar hyperbolic equation. *SIAM J. Numer. Anal.*, 28:133–140, 1991.
 - [42] M. Portelheiro and A. Tzvaras. Hydrodynamic limits for kinetic equations and the diffusive approximation of radiative transport for acoustic waves. *Trans. Amer. Math. Soc.*, 359:529–565, 2007.
 - [43] W. Reed and T. Hill. Triangular mesh methods for the neutron transport equation. Technical

- report, Los Alamos National Laboratory, Los Alamos, NM, 1973.
- [44] L. Reggiani. *Hot-electron transport in semiconductors*, volume 58 of *Topics in Applied Physics*. Springer, Berlin, 1985.
 - [45] G. R. Richter. An optimal-order error estimate for the discontinuous Galerkin method. *Math. Comp.*, 50:75–88, 1988.
 - [46] C.-W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock-capturing schemes. *J. Comput. Phys.*, 77:439–471, 1988.
 - [47] R. M. Strain and Y. Guo. Almost exponential decay near Maxwellian. *Comm. Partial Differential Equations*, 31(1-3):417–429, 2006.
 - [48] K. Tomizawa. *Numerical simulation of sub micron semiconductor devices*. Artech House, Boston, 1993.
 - [49] C. Villani. A review of mathematical topics in collisional kinetic theory handbook of mathematical fluid dynamics. I:71–74, 2002.
 - [50] C. Villani. Hypocoercive diffusion operators in Hörmander form. *Mathematical Models and Methods in Applied Sciences*, 2006.
 - [51] X. Zhang and C.-W. Shu. A genuinely high order total variation diminishing scheme for one-dimensional scalar conservation laws. *SIAM Journal on Numerical Analysis*, 48(2):772–795, 2010.
 - [52] X. Zhang and C.-W. Shu. On maximum-principle-satisfying high order schemes for scalar conservation laws. *J. Comput. Phys.*, 229:3091–3120, 2010.
 - [53] X. Zhang and C.-W. Shu. On positivity preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes. *Journal of Computational Physics*, 229:8918–8934, 2010.
 - [54] X. Zhang, Y. Xia, and C.-W. Shu. Maximum-principle-satisfying and positivity-preserving high order discontinuous Galerkin schemes for conservation laws on triangular meshes. *Journal of Scientific Computing*. preprint submitted.
 - [55] X. Zhang, Y. Xing, and C.-W. Shu. Positivity preserving high order well balanced discontinuous Galerkin methods for the shallow water equations. *Advances in Water Resources*. to appear.

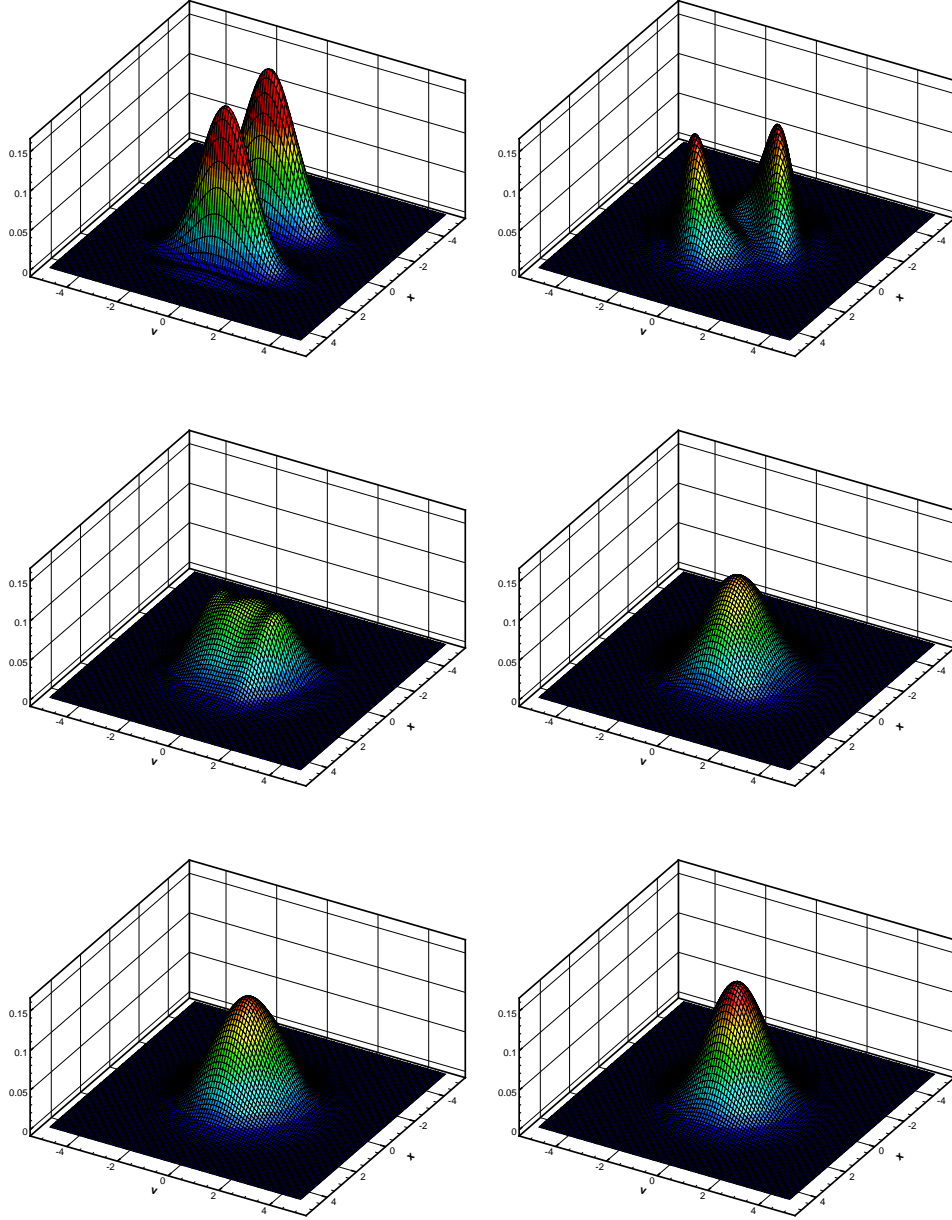


FIG. 6.3. The evolution of pdf to equilibrium. Top left: $t=0$, top right: $t=1$, middle left: $t=2$, middle right: $t=3$, bottom left: $t=4$, bottom right: $t=6$. Positivity-preserving DG scheme computed on a 50×50 mesh with piecewise quadratic polynomials and third order Runge-Kutta time stepping.