

# A successive penalty–based Asymptotic-Preserving scheme for kinetic equations <sup>\*</sup>

Bokai Yan<sup>†</sup>                      Shi Jin<sup>‡</sup>

September 30, 2012

## Abstract

We propose an asymptotic-preserving (AP) scheme for kinetic equations that is efficient also in the hydrodynamic regimes. This scheme is based on the BGK-penalty method introduced by Filbet-Jin [14], but uses the penalization successively to achieve the desired asymptotic property. This method possesses a stronger AP property than the original method of Filbet-Jin, with the additional feature of being also positivity preserving when applied on the Boltzmann equation. It is also general enough to be applicable to several important classes of kinetic equations, including the Boltzmann equation and the Landau equation. Numerical experiments verify these properties.

## 1 Introduction

In the study of rarefied gas or plasma physics, the distribution function  $f(t, x, v)$  is usually used to describe the density of the particles at time  $t$  and position  $x$ , with velocity  $v$ . This distribution can be modeled by the kinetic equation [6],

$$\frac{\partial f}{\partial t} + v \cdot \nabla_x f = \frac{1}{\varepsilon} Q(f), \quad (1)$$

For molecules with primarily binary short range collisions,  $Q(f)$  is given by the Boltzmann collision operator

$$Q(f) = \iint_{\mathbb{R}^N \times \mathbb{S}^{N-1}} B(|v - v_*|, \cos \theta) (f' f'_* - f f_*) dv_* d\sigma. \quad (2)$$

Here we use the shorthanded notations  $f = f(v)$ ,  $f_* = f(v_*)$ ,  $f' = f(v')$  and  $f'_* = f(v'_*)$ . The post-collision velocities can be computed by

$$v' = v - \frac{1}{2} \left( (v - v_*) - |v - v_*| \sigma \right), \quad v'_* = v - \frac{1}{2} \left( (v - v_*) + |v - v_*| \sigma \right),$$

where  $\sigma$  is a unit vector varying in the sphere  $\mathbb{S}^{N-1}$ .

---

<sup>\*</sup>This work was partially supported by NSF grants DMS-0608720, DMS-1114546 and NSF FRG grant DMS-0757285. SJ is also supported by a Vilas Associate Award from the University of Wisconsin-Madison.

<sup>†</sup>Department of Mathematics, University of Wisconsin-Madison, Madison, WI 53706, USA (yan@math.wisc.edu)

<sup>‡</sup>Department of Mathematics, Institute of Natural Sciences, and Ministry of Education Key Laboratory of Scientific and Engineering Computing, Shanghai Jiao Tong University, Shanghai 200240, China; and Department of Mathematics, University of Wisconsin-Madison, Madison, WI 53706, USA (jin@math.wisc.edu)

We focus on the *Variable Hard Sphere* model in this work,

$$B(|u|, \cos \theta) = C|u|^r.$$

For charged particles in plasma physics, where the long range interaction dominates,  $Q(f)$  is given by the Landau collision operator [26, 27]

$$Q(f) = \nabla_v \cdot \int_{\mathbb{R}^{N_v}} A(v - v_*) (f(v_*) \nabla_v f(v) - f(v) \nabla_{v_*} f(v_*)) dv_*, \quad (3)$$

where the semi-positive definite matrix  $A(z)$  is given by

$$A(z) = \Psi(z) \left( I - \frac{z \otimes z}{|z|^2} \right), \quad \Psi(z) = |z|^{\gamma+2}. \quad (4)$$

The parameter  $\gamma$  is determined by the type of interaction between particles. In the case of inverse power law relationship, that is, when two particles at distance  $r$  interact with a force proportional to  $1/r^s$ ,  $\gamma = \frac{s-5}{s-1}$ . For example, in the cases of the Maxwell molecules  $\gamma = 0$  (corresponding to  $s = 5$ ) and for the Coulomb potential  $\gamma = -3$  (corresponding to  $s = 2$ ). The Landau equation is derived as a limit of the Boltzmann equation when all the collisions become grazing [1, 7, 9, 18, 32]. We refer to [25] and references therein for details. In this article we will always take  $\gamma = -3$ .

Both operators (2) and (3) satisfy some important properties:

- Conservation of mass, momentum and energy:

$$\int \phi Q(f) dv = 0, \quad \text{with } \phi = 1, v, \frac{|v|^2}{2};$$

- Entropy dissipation:

$$\frac{d}{dt} \int f \log f dv = \int Q(f) \log f dv \leq 0$$

with equality holds if and only if  $f = M$ ;

- Well balancedness:

$$Q(f) = 0 \Leftrightarrow f = M.$$

Here the Maxwellian  $M$  is the local equilibrium,

$$M = \frac{\rho}{(2\pi T)^{N/2}} \exp\left(-\frac{(v-u)^2}{2T}\right), \quad (5)$$

with density  $\rho$ , macroscopic velocity  $u$  and temperature  $T$  defined by

$$\rho = \int f dv, \quad \rho u = \int v f dv, \quad \rho T = \frac{1}{N} \int |v-u|^2 f dv.$$

The Knudsen number  $\varepsilon$  in (1) is the ratio between the mean free path and the typical physical length scale.

As  $\varepsilon \rightarrow 0$ , the moments of solution to (1) can be approximated by the macroscopic compressible Euler equations [2], [5],

$$\begin{cases} \partial_t \rho + \nabla_x \cdot \rho u = 0 \\ \partial_t(\rho u) + \nabla_x \cdot (\rho u \otimes u + pI) = 0 \\ \partial_t E + \nabla_x \cdot ((E+p)u) = 0 \end{cases} \quad (6)$$

with total energy  $E$

$$E = \frac{1}{2} \rho u^2 + \frac{N}{2} \rho T = \int \frac{|v|^2}{2} f dv$$

and pressure  $p$  given by the constitutive relation to close the system (6)

$$p = \rho T.$$

Since the kinetic equation (1) approaches the Euler equations (6) asymptotically as  $\varepsilon \rightarrow 0$ , it is a natural request that a good scheme designed for (1) can capture the asymptotic limit (6) as  $\varepsilon \rightarrow 0$ , with time step and mesh sizes in space and velocity spaces fixed. A scheme with this property is called Asymptotic Preserving (AP) [22]. Such schemes are able to remove the stiffness for small  $\varepsilon$ , and can capture the macroscopic hydrodynamic behavior without numerically resolving the small  $\varepsilon$ . There have been active recent research activities in developing AP schemes for Boltzmann equation, see [11, 12, 3, 10, 28]. We refer to a recent review [23] on AP schemes for kinetic and hyperbolic equations.

More specifically, for Boltzmann type equation, let  $f^n$  be the numerical solution approximating  $f(t^n)$ . Then the AP property is equivalent to require

$$f^n - M^n = \mathcal{O}(\varepsilon), \quad \text{for any } n \geq 1, \quad (7)$$

for any initial data, equilibrium or non-equilibrium. As in [23], we call this result a *strong AP* property.

One of the main challenges to the development of AP schemes for the kinetic equation (1) is the implicit collision term, if the time step is required to be larger than  $\varepsilon$ . The collision operator  $Q(f)$  is typically nonlinear, nonlocal and high dimensional, thus its numerical inversion is computationally difficult and expensive. Recently Filbet and Jin [14] introduced a BGK-penalty method for (1) with the Boltzmann operator (2) that overcomes this difficulty. The idea was to penalize  $Q(f)$  by the BGK-operator

$$P(f) = M - f, \quad (8)$$

which can be inverted easily, and treat  $Q(f)$  *explicitly*. This results a scheme that has the *relaxed AP* property in the sense that for any  $\varepsilon > 0$ , there exists an integer  $N > 0$ ,

$$f^n - M^n = \mathcal{O}(\varepsilon), \quad \text{for any } n \geq N. \quad (9)$$

This means that the AP property is satisfied after an initial transient time.

Later the authors [25] extended this result to the nonlinear Landau equation (1) with (3), based on a Fokker-Planck penalization. A similar relaxed AP property was obtained. This method was also extended to the quantum cases [13, 20] and multispecies case [24], and also in diffusive limit of linear transport equation [8]. A rigorous analysis on the application of this method to hyperbolic system was given in [17].

The BGK-penalty method can be implemented differently. In the work of Dimarco and Pareschi [10] for the Boltzmann equation, a time-splitting approach was introduced, so the convection is solved in a separate step from the collision step. The collision step, using the fact that the local Maxwellian is invariant for the space homogeneous Boltzmann equation, can be solved using the exponential Runge-Kutta method. Their method is positivity-preserving and has the *exponential AP* property, in the sense that there exists some constant  $c > 0$ , such that for any initial data,

$$f^n - M^n = \mathcal{O}(e^{-\frac{c\Delta t}{\varepsilon}}), \quad \text{for any } n \geq 1. \quad (10)$$

This method is also generalized to the diffusive limit [4]. However this method has not been extended to the Landau equation (1)(3), or other more general collision operators, which cannot take advantage of the special property of the BKG operator in the space homogeneous case.

The goal of this work is to improve the Filbet-Jin method in two aspects: 1) positivity preserving and 2) a strong AP property (7). This new method is based on *successive* penalizations, namely the penalty operators will be utilized more than once in each time step. This gives a method in between that of Filbet-Jin and Dimarco-Pareschi, and that combines the advantages of both methods in

terms of positivity, asymptotic-preserving and generality. With the positivity this formulation is also suitable for a Monte-Carlo simulation [10].

This paper is organized as follows. We briefly review the penalty method of Filbet-Jin and the exponential method of Dimarco-Pareschi in section 2, and introduce a positivity-preserving improvement for the Filbet-Jin method. In section 3 we introduce the new successive penalty method in the first and second order formulations, and study its asymptotic property. Finally numerical experiments are carried out in section 4 for both the Boltzmann and Landau equations to study the properties of the new method. The paper is concluded in section 5.

## 2 Penalty based methods

### 2.1 The Filbet-Jin method

#### 2.1.1 For the Boltzmann equation

First we briefly review the work of Filbet and Jin [14] for the Boltzmann equation (1)(2). The idea in Filbet-Jin's method is to penalize  $Q(f)$  by another operator  $P(f)$ ,

$$\frac{\partial f}{\partial t} + v \cdot \nabla_x f = \underbrace{\frac{1}{\varepsilon}(Q(f) - \beta P(f))}_{\text{less stiff}} + \underbrace{\frac{1}{\varepsilon}\beta P(f)}_{\text{stiff}}, \quad (11)$$

then the less stiff term can be solved explicitly and the new stiff term is solved implicitly. A scheme with first order accuracy in time reads

$$\frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^n = \frac{1}{\varepsilon} (Q(f^n) - \beta P(f^n) + \beta P(f^{n+1})). \quad (12)$$

A second order scheme is obtained by

$$\begin{cases} \frac{f^* - f^n}{\Delta t/2} + v \cdot \nabla_x f^n = \frac{Q(f^n) - \beta^n P(f^n)}{\varepsilon} + \frac{\beta^n P(f^*)}{\varepsilon}, \\ \frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^* = \frac{Q(f^*) - \beta^* P(f^*)}{\varepsilon} + \frac{\beta^*(P(f^n) + P(f^{n+1}))}{2\varepsilon}. \end{cases} \quad (13)$$

One wants  $P(f)$  to be easy to invert while at the same time to preserve the good properties of  $Q(f)$ . A good choice used by Filbet-Jin is the BGK operator (8). Then scheme (12) reads

$$\frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^n = \frac{1}{\varepsilon} (Q(f^n) - \beta(M^n - f^n) + \beta(M^{n+1} - f^{n+1})). \quad (14)$$

$M^{n+1}$  can be solved explicitly first thanks to the fact that the right side of (14) preserves density, momentum and energy. Multiplying  $\phi = 1, v, \frac{|v|^2}{2}$  to (14) and integrating over velocity space, one obtains

$$\int \phi \left( \frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^n \right) dv = 0.$$

Then the moments at  $t^{n+1}$  can be derived explicitly,

$$(\rho, \rho u, E)^{n+1} = \int \phi (f^n - \Delta t v \cdot \nabla_x f^n) dv. \quad (15)$$

and  $M^{n+1}$  is obtained. Then  $f^{n+1}$  can be solved,

$$f^{n+1} = M^{n+1} + \frac{1}{1 + \frac{\beta \Delta t}{\varepsilon}} \left( f^n - \Delta t v \cdot \nabla_x f^n - M^{n+1} + \frac{\Delta t}{\varepsilon} (Q(f^n) - \beta(M^n - f^n)) \right). \quad (16)$$

Therefore the implicit scheme (12) can be solved *explicitly*. The implementation for (13) is similar. The stability condition for (14) can be derived,

$$\beta > \frac{1}{2} \|\nabla Q(f)\|_\infty, \quad (17)$$

where  $\nabla Q(f)$  is the Frechet derivative of  $Q$  around the corresponding Maxwellian.

At last we give the weak-AP property proved in [14].

**Theorem 2.1** *Consider the numerical solution given by (14). Then*

1. *If  $\varepsilon \rightarrow 0$  and  $f^n = M^n + \mathcal{O}(\varepsilon)$ , then  $f^{n+1} = M^{n+1} + \mathcal{O}(\varepsilon)$ . Thus, when  $\varepsilon \rightarrow 0$ , the (moments of the) scheme becomes a consistent discretization of the Euler system (6).*
2. *Assume  $\varepsilon \ll 1$  and  $f^n = M^n + \mathcal{O}(\varepsilon)$ . If there exists a constant  $C > 0$  such that*

$$\left\| \frac{f^{n+1} - f^n}{\Delta t} \right\| + \left\| \frac{U^{n+1} - U^n}{\Delta t} \right\| \leq C,$$

*then the scheme asymptotically becomes a first order in time approximation of the compressible Navier-Stokes equations.*

### 2.1.2 A positivity-preserving improvement

For numerical purpose, we assume that  $f$  has a compact support in  $\Omega_V = [-v_{\max}, v_{\max}]^N \in \mathbb{R}^N$  in  $v$  direction. The computation in this article is always performed on  $\Omega_V$ . However the results can be extended to any other compact domain.

One question unsolved in the Filbet-Jin paper is the positivity of the scheme. More specifically, when initial data  $f^I$  is nonnegative over  $\mathbb{R}^N \times \Omega_V$ , one hopes the distribution  $f$  is always nonnegative during the time evolution.

(14) can be positive preserving after a small correction. The key idea is that the two  $\beta$ 's in (14) do not have to have the same value. A difference of  $\mathcal{O}(\Delta t)$  is permitted to keep the first order convergence. A simple calculation shows that the scheme is positive if one puts a little more weight on the second  $\beta$ . Besides, all the other good properties of (14), like AP and stability, remain valid.

Note that the Boltzmann operator (2) can be split to a gain term and a loss term:

$$Q(f) = Q^+(f) - fQ^-(f), \quad (18)$$

with

$$\begin{aligned} Q^+(f) &= \iint_{\mathbb{R}^N \times \mathbb{S}^{N-1}} C|v - v_*|^r f' f'_* dv_* d\sigma, \\ Q^-(f) &= \iint_{\mathbb{R}^N \times \mathbb{S}^{N-1}} C|v - v_*|^r f_* dv_* d\sigma. \end{aligned} \quad (19)$$

Consider the first order scheme (12) with the Boltzmann collision operator  $Q(f) = Q^+(f) - fQ^-(f)$  and the BGK operator  $\beta P(f) = \beta(M - f)$ . We choose different  $\beta$  for each  $P(f)$ ,

$$\frac{f^{n+1} - f^n}{\Delta t} + v \cdot D_x f^n = \frac{1}{\varepsilon} (Q(f^n) - \beta^n(M^n - f^n) + \beta^n(1 + \Delta t \kappa^n)(M^{n+1} - f^{n+1})), \quad (20)$$

where  $D_x$  is some positive preserving discretization of  $\nabla_x$  (for example the upwind scheme).  $\beta^n$  and  $\kappa^n$  are  $x$ -dependent only and given by

$$\beta^n = \max_v Q^-(f^n), \quad (21)$$

$$\kappa^n = \max\left\{ \max_v \frac{-(M^{n+1} - M^n)}{\Delta t M^{n+1}}, 0 \right\}, \quad (22)$$

where  $v \in \Omega_V$  is bounded.

**Lemma 2.2** *The time discrete scheme (20)(21)(22) is well defined and first order in time. Besides, if the CFL condition  $v_{\max}\Delta t \leq \Delta x$  is satisfied, then,*

1. *If the initial data is close to the local Maxwellian  $f^I = M^I + \mathcal{O}(\varepsilon)$ , then the scheme is asymptotic preserving.*
2. *If the initial data is nonnegative, then  $f^n$  remains positive, for any  $n \geq 1$ .*

**Proof.** One can easily find an upper bound for  $Q^-(f^n)$ ,

$$Q^-(f^n) = \int_{\Omega_V \times \mathbb{S}^{N-1}} C|v - v_*|^r f_* dv_* d\sigma \leq C|v_{\max}|^r \int_{\Omega_V} f_* dv_* = C\rho,$$

which gives a well defined  $\beta^n$  by (21).

Next, one can compute the term  $M^{n+1}$  without  $\kappa^n$ , due to the fact that  $\beta^n$  and  $\kappa^n$  are not  $v$  dependent. More specifically, this can be done by multiplying (20) with  $\phi = 1, v, |v|^2/2$  and integrating with respect to  $v$ , which gives exactly (15). Then  $M^{n+1}$  is defined and  $\kappa^n$  can be found by (22).

To show that the scheme is first order in time, one only needs  $\kappa^n = \mathcal{O}(1)$  which is true since

$$\frac{M^n - M^{n+1}}{\Delta t M^{n+1}} \approx \frac{\log M^n - \log M^{n+1}}{\Delta t} = \mathcal{O}(1).$$

A more precise estimate on  $\kappa$  is given by the following remark.

**Remark 2.3** *One might expect that the value of  $\kappa$  could be very large since its definition has a term  $M^{n+1}$  in the denominator, which is close to 0 near the artificial boundary  $\{|v| = v_{\max}\} \subset V$ . However the value of  $M^n/M^{n+1}$  is fairly close to  $\rho^n(T^{n+1})^{N/2}/(\rho^{n+1}(T^n)^{N/2})$ . It does not “blow up” near the artificial boundary. We leave a detailed computation in the appendix.*

The AP property is part of Theorem 2.1.

Next we show the positivity, when the CFL condition is satisfied.

Suppose  $f^n$  is nonnegative.  $f^{n+1}$  can be solved from (20)

$$\left(1 + \frac{\Delta t \beta^n (1 + \Delta t \kappa^n)}{\varepsilon}\right) f^{n+1} = (f^n - \Delta t v \cdot D_x f^n) + \frac{\Delta t}{\varepsilon} (Q^+(f^n) + (\beta^n - Q^-(f^n))f^n + (\beta^n(1 + \Delta t \kappa^n)M^{n+1} - \beta^n M^n).$$

The transport term  $(f^n - \Delta t v \cdot D_x f^n)$  is positive if the CFL condition is satisfied. The term  $Q^+(f^n)$  is positive by its definition. To get a positive  $f^{n+1}$ , one also needs

$$\beta^n - Q^-(f^n) \geq 0,$$

$$\kappa^n + \frac{M^n - M^{n+1}}{\Delta t M^{n+1}} \geq 0.$$

Clearly these conditions are satisfied if one chooses  $\beta^n$  and  $\kappa^n$  as in (21)(22).

**Remark 2.4** *From the proof of Theorem 2.2, a sufficient condition for (20) to be positive preserving is, for any  $v$ ,*

$$\begin{aligned} \beta^n &\geq Q^-(f^n), \\ \kappa^n &\geq -\frac{M^n - M^{n+1}}{\Delta t M^{n+1}}. \end{aligned} \tag{23}$$

However, larger  $\beta$  and  $\kappa$  can reduce the accuracy of the scheme. A simple numerical analysis shows that the local truncation error is given by

$$f((n+1)\Delta t) - f^{n+1} = \frac{\Delta t^2}{2} \partial_{tt} f^n + \beta^n \frac{\Delta t^2}{\varepsilon} \left( \kappa^n f^n - v \cdot \nabla_x f^n + \frac{Q^n}{\varepsilon} - \left( \frac{M^{n+1} - M^n}{\Delta t} + \kappa^n M^{n+1} \right) \right) + \text{low order terms.}$$

Therefore (21)(22) give the best choices.

**Remark 2.5** The positive preserving technique cannot be applied directly to the second order scheme (13). The main reason is that, the IMEX scheme used in (13) is not positive preserving, even if the penalization technique is not used and the Boltzmann collision can be solved fully implicitly. The transport part in the second equation of (13) is discretized at time  $t^*$ , instead of  $t^n$ , which introduces uncontrolled negative parts when plugging in the  $f^*$  obtained from the first equation of (13).

### 2.1.3 For the Landau equation

The Filbet-Jin method was extended to the Landau equation (1)(3) in [25]. The BGK operator (8) is not a suitable choice for penalization for this equation, since the diffusive nature of the Landau operator (3) introduces extra stiffness. Instead the Fokker-Planck operator was used:

$$P_{FP}(f) = P_{FP}^M f = \nabla_v \cdot \left( M \nabla_v \left( \frac{f}{M} \right) \right). \quad (24)$$

The first order scheme reads

$$\frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^n = \frac{1}{\varepsilon} (Q(f^n) - \beta P^n f^n + \beta P^{n+1} f^{n+1}) \quad (25)$$

where  $P^n f^n = P_{FP}^{M^n} f^n$  is the FP operator (24) and  $\beta$  is given by

$$\beta = \beta_0 \max_v \lambda(D_A(f)). \quad (26)$$

Here  $\beta_0$  is a constant satisfying  $\beta_0 > \frac{1}{2}$ . A good choice is  $\beta_0 = 1$ .  $\lambda(D_A)$  is the spectral radius of the positive symmetric matrix  $D_A$ , with  $D_A(f)$  defined by

$$D_A(f) = \int A(v - v_*) f_* dv_*. \quad (27)$$

A second order implicit-explicit (IMEX) type scheme reads

$$\begin{cases} \frac{f^* - f^n}{\Delta t/2} + v \cdot \nabla_x f^n = \frac{Q(f^n) - \beta P^n f^n}{\varepsilon} + \frac{\beta P^* f^*}{\varepsilon}, \\ \frac{f^{n+1} - f^n}{\Delta t} + v \cdot \nabla_x f^* = \frac{Q(f^*) - \beta^* P^* f^*}{\varepsilon} + \frac{\beta^* P^n f^n + \beta^* P^{n+1} f^{n+1}}{2\varepsilon}. \end{cases} \quad (28)$$

with  $P(f)$  the FP operator (24). Suggested by numerical experiments, one can take

$$\begin{aligned} \beta &= \beta_0 \max_{v,\lambda} \lambda(D_A(f)), \\ \beta^* &= \beta_0 \max_{v,\lambda} \lambda(D_A(f^*)). \end{aligned} \quad (29)$$

Again the constant coefficient satisfies  $\beta_0 > \frac{1}{2}$ . A good choice is  $\beta_0 = (2 + \sqrt{2})$ .

An efficient method to invert the FP operator  $P_{FP}$  was also introduced in [25].

**Remark 2.6** Unfortunately we cannot derive a positive preserving method for Landau equation at this time. The technique introduced for Boltzmann equation in Section 2.1.2 cannot be applied here. In fact, to the authors' best knowledge, there are no conservative methods yet which implicitly solve the Landau equation with the positive preserving property. We refer to [29] and references therein for some implicit Landau solvers.

## 2.2 The Dimarco-Pareschi method for the Boltzmann equation

Utilizing on this BGK penalization, Dimarco and Pareschi introduced a class of exponential Runge-Kutta methods in [10] for the Boltzmann equation, which are exponentially AP in the sense of (10). The starting point is to split the Boltzmann equation (1) into a relaxation step

$$\frac{\partial f}{\partial t} = \frac{1}{\varepsilon} Q(f), \quad (30)$$

and a transport step,

$$\frac{\partial f}{\partial t} + v \cdot \nabla_x f = 0. \quad (31)$$

Ignoring the convection operator in (11), then (11) and (8) can be written as

$$\frac{\partial f}{\partial t} = \frac{1}{\varepsilon} \left( \tilde{Q}(f) - \beta M \right) + \frac{\beta}{\varepsilon} (M - f), \quad (32)$$

where

$$\tilde{Q}(f) = Q(f) + \beta f,$$

with some constant  $\beta$ .

Noting that the macroscopic quantities  $\rho$ ,  $u$  and  $T$  (hence  $M$ ) are not changed in this step. (32) can be reformulated as

$$\frac{\partial (f - M) e^{\beta t / \varepsilon}}{\partial t} = \frac{1}{\varepsilon} \left( \tilde{Q}(f) - \beta M \right) e^{\beta t / \varepsilon}. \quad (33)$$

A class of explicit exponential Runge-Kutta schemes can be obtained. For example, one can apply the explicit Euler method to this system

$$\frac{(f^* - M^*) e^{\beta(t^n + \Delta t) / \varepsilon} - (f^n - M^n) e^{\beta t^n / \varepsilon}}{\Delta t} = \frac{1}{\varepsilon} \left( \tilde{Q}(f^n) - \beta M^n \right) e^{\beta t^n / \varepsilon}.$$

Since  $M^* = M^n$ , one obtains,

$$f^* = e^{-\beta \Delta t / \varepsilon} f^n + \frac{\beta \Delta t}{\varepsilon} e^{-\beta \Delta t / \varepsilon} \frac{\tilde{Q}(f^n)}{\beta} + \left( 1 - \left( 1 + \frac{\Delta t}{\varepsilon} \beta \right) e^{-\beta \Delta t / \varepsilon} \right) M^n. \quad (34)$$

Then the transport step (31) can be solved by an explicit scheme, for example the upwind method.

As  $\varepsilon \rightarrow 0$ , one has  $f^* = M^n$ . Then the moments of the transport step give a kinetic scheme for the Euler system (6). One obtains an exponentially AP scheme, in the sense of (10) with the constant  $c = \beta$ .

The positivity of  $f^*$  is guaranteed as long as  $\tilde{Q}(f^n)$  is positive, which holds under the condition

$$\beta^n \geq Q^-(f^n). \quad (35)$$

This is exactly the first equation in (23).

A remarkable feature is that, (34) solves  $f^*$  as a convex combination of positive functions  $f^n$ ,  $\tilde{Q}(f^n)$  and  $M^n$ . Hence the Monte Carlo technique can be applied based on this formulation (see [10]).



Higher order schemes can be derived by applying high order temporal operator splitting on (1), high order Runge-Kutta method on the system (33) and high order methods on the transport equation (31). See [10] for details.

The extension of the Dimarco-Pareschi method to the Landau equation (1)(3) is not easy, since the exact solution of Fokker-Planck operator  $P$  is not easy to find. The Filbet-Jin method requires the (implicit) numerical solution, which is relatively easier than the Dimarco-Pareschi method. However, as discussed before, only a relaxed AP property is obtained for the Filbet-Jin method.

### 3 A successive penalty method

Let us think about these two penalty methods in a different way.

#### 3.1 A toy model

Consider the toy model,

$$\frac{df}{dt} = -\frac{1}{\varepsilon}f. \quad (36)$$

We can apply the Filbet-Jin method and the Dimarco-Pareschi method on this equation. Both methods start with the reformulation

$$\frac{df}{dt} = -\frac{1-\beta}{\varepsilon}f - \frac{\beta}{\varepsilon}f.$$

After a time splitting, one obtains

$$\frac{df}{dt} = -\frac{1-\beta}{\varepsilon}f, \quad (37)$$

$$\frac{df}{dt} = -\frac{\beta}{\varepsilon}f. \quad (38)$$

(37) is a non-stiff (or less stiff) part, hence solved explicitly,

$$\frac{f^* - f^n}{\Delta t} = -\frac{1-\beta}{\varepsilon}f^n.$$

The difference of the two methods lies in how to solve the stiff part (38). The Filbet-Jin method solves this step *implicitly*,

$$\frac{f^{n+1} - f^*}{\Delta t} = -\frac{\beta}{\varepsilon}f^{n+1}.$$

Therefore

$$f^{n+1} = \frac{1 + \frac{\beta-1}{\varepsilon}\Delta t}{1 + \frac{\beta}{\varepsilon}\Delta t} f^n. \quad (39)$$

The Dimarco-Pareschi method solves this step *exactly*,

$$f^{n+1} = e^{-\frac{\beta}{\varepsilon}\Delta t} f^*.$$

Therefore

$$f^{n+1} = \frac{1 + \frac{\beta-1}{\varepsilon}\Delta t}{e^{\frac{\beta}{\varepsilon}\Delta t}} f^n. \quad (40)$$

It is natural to design a method which solves the stiff part (38) in a different way. One can divide the time interval  $[t^n, t^{n+1}]$  into  $k$  subintervals, and apply the implicit Euler method in each

subinterval, i.e.

$$\left\{ \begin{array}{l} \frac{f^{n+1,1} - f^*}{\Delta t/k} = -\frac{\beta}{\varepsilon} f^{n+1,1}, \\ \frac{f^{n+1,2} - f^{n+1,1}}{\Delta t/k} = -\frac{\beta}{\varepsilon} f^{n+1,2}, \\ \dots \\ \frac{f^{n+1} - f^{n+1,k-1}}{\Delta t/k} = -\frac{\beta}{\varepsilon} f^{n+1}. \end{array} \right.$$

Hence

$$f^{n+1} = \left(1 + \frac{\beta \Delta t}{\varepsilon k}\right)^{-k} f^*.$$

Therefore

$$f^{n+1} = \frac{1 + \frac{\beta-1}{\varepsilon} \Delta t}{\left(1 + \frac{\beta \Delta t}{\varepsilon k}\right)^k} f^n. \quad (41)$$

Noting that

$$\left(1 + \frac{\beta \Delta t}{\varepsilon k}\right)^k \geq 1 + \frac{\beta}{\varepsilon} \Delta t,$$

this method is unconditionally stable if  $\beta \geq \frac{1}{2}$ . The positivity is preserved under a stronger condition  $\beta \geq 1$ .

When  $k = 1$ , this gives the Filbet-Jin method (39). When  $k \rightarrow \infty$ , this gives the Dimarco-Pareschi method (40). Here we take  $k = 2$ , which gives an intermediate method between these two methods. In this case one obtains the strong AP property

$$f^n = \mathcal{O}(\varepsilon), \quad \text{for any } n \geq 1.$$

We call this the *successive penalty* method, due to the fact that the implicit part is solved in two (or more) successive steps.

### 3.2 A successive penalty method for kinetic equations

We can apply this idea to kinetic equation (1) with a penalization operator  $P$ .

The Dimarco-Pareschi method applies an operator splitting between the relaxation step and the transport step, while the Filbet-Jin method is based on an unsplit version. It turns out that whether to apply this operator splitting plays an important role.

#### The split version

The operator splitting between the relaxation step and the transport step is necessary for the Dimarco-Pareschi method since the key idea in their method is that the BGK operator  $P$  can be solved exactly when the Maxwellian  $M$  is time independent. With this splitting, we give the following successive penalty method,

$$\left\{ \begin{array}{l} \frac{f^* - f^n}{\Delta t} = \frac{Q(f^n) - \beta P(f^n)}{\varepsilon} + \frac{(1-\alpha)\beta P(f^*)}{\varepsilon}, \\ \frac{f^{**} - f^*}{\Delta t} = \frac{\alpha\beta P(f^{**})}{\varepsilon}, \\ \frac{f^{n+1} - f^{**}}{\Delta t} + v \cdot \nabla_x f^{**} = 0, \end{array} \right. \quad (42)$$

with a constant  $\alpha \in (0, 1)$ . One can simply choose  $\alpha = \frac{1}{2}$ , as what we do for the toy model.

This can be seen as an approximation of the Dimarco-Pareschi method, with easier extension to more complicated problems. (42) can be applied to both the Boltzmann equation and the Landau equation, with the penalization  $P$  to be the BGK operator (8) or the Fokker-Planck operator (24), and the penalization weight  $\beta$  given by (17) or (26), respectively.

It is easy to show that the strong AP property (7) is satisfied.

In the case of the Boltzmann equation, i.e.,  $Q(f)$  is the Boltzmann operator (2) and  $P(f)$  is the BGK operator, the relaxation step gives,

$$f^{**} = Bf^n + \frac{\beta\Delta t}{\varepsilon} B \frac{\tilde{Q}(f^n)}{\beta} + \left(1 - B - \frac{\beta\Delta t}{\varepsilon} B\right) M^n, \quad (43)$$

where

$$B = \frac{1}{\left(1 + \frac{\alpha\beta\Delta t}{\varepsilon}\right) \left(1 + \frac{(1-\alpha)\beta\Delta t}{\varepsilon}\right)}.$$

Noting that  $\tilde{Q}(f^n) = Q(f^n) + \beta f^n$  is non-negative under the condition (35) and  $\left(1 + \frac{\beta\Delta t}{\varepsilon}\right) B \leq 1$ , (43) also solves  $f^{**}$  as a convex combination of positive functions  $f^n$ ,  $\frac{\tilde{Q}(f^n)}{\beta}$  and  $M^n$ , as in the Dimarco-Pareschi method (34). Hence  $f^{**}$  is non-negative and the Monte Carlo technique can be applied.

### The nonsplit version

Following the Filbet-Jin method, we can give the successive penalty method without operator splitting:

$$\begin{cases} \frac{f^* - f^n}{\Delta t} + v \cdot \nabla_x f^n = \frac{Q(f^n) - \beta P(f^n)}{\varepsilon} + \frac{(1 - \alpha^n)\beta P(f^*)}{\varepsilon}, \\ \frac{f^{n+1} - f^*}{\Delta t} = \frac{\alpha^n \beta P(f^{n+1})}{\varepsilon}, \end{cases} \quad (44)$$

where the time dependent  $\alpha^n \in (0, 1)$  will be specified later.

Note that the solution is given by

$$f^{n+1} = \left(1 - \alpha \frac{\beta\Delta t}{\varepsilon} P\right)^{-1} \left(1 - (1 - \alpha) \frac{\beta\Delta t}{\varepsilon} P\right)^{-1} \left(f^n - \Delta t v \cdot D_x f^n + \frac{\Delta t}{\varepsilon} (Q(f^n) - \beta P(f^n))\right).$$

If one takes a constant  $\alpha$ , with initial data  $f^0 = M^0 + \mathcal{O}(\varepsilon)$ , one would have a much stronger AP property

$$f^n = M^n + \mathcal{O}(\varepsilon^2). \quad (45)$$

This is between the strong AP property (7) and the exponential AP property (10). To derive the typical strong AP property (7), one can choose a time dependent  $\alpha$  which is  $\mathcal{O}(\varepsilon)$  when  $f$  is close to the equilibrium  $M$ . In practice the following choice works well

$$\alpha^n = \alpha^n(x) = \min \left\{ \frac{\|f^n - M^n\|}{\Delta t}, \frac{1}{2} \right\}, \quad (46)$$

where the norm  $\|\cdot\|$  is taken over the velocity space.

With this choice one can show that the strong AP property (7) is satisfied.

**Remark 3.1** *In the case of the Boltzmann equation, i.e.,  $Q(f)$  is the Boltzmann operator (2) and  $P(f)$  is the BGK operator, (44) gives,*

$$\begin{cases} \frac{f^* - f^n}{\Delta t} + v \cdot \nabla_x f^n = \frac{Q(f^n) - \beta(M^n - f^n)}{\varepsilon} + (1 - \alpha) \frac{\beta(M^{n+1} - f^*)}{\varepsilon}, \\ \frac{f^{n+1} - f^*}{\Delta t} = \alpha \frac{\beta(M^{n+1} - f^{n+1})}{\varepsilon}. \end{cases} \quad (47)$$

This is solved in a similar way as in the Filbet-Jin method. The result is

$$f^{n+1} = M^{n+1} + \frac{1}{\left(1 + \frac{\alpha\beta\Delta t}{\varepsilon}\right) \left(1 + \frac{(1-\alpha)\beta\Delta t}{\varepsilon}\right)} \left( f^n - \Delta t v \cdot D_x f^n - M^{n+1} + \frac{\Delta t}{\varepsilon} (Q(f^n) - \beta(M^n - f^n)) \right). \quad (48)$$

Compared with the Filbet-Jin method (16), we simply change the bottom  $\left(1 + \frac{\beta\Delta t}{\varepsilon}\right)$  to a larger number  $\left(1 + \frac{\alpha\beta\Delta t}{\varepsilon}\right) \left(1 + \frac{(1-\alpha)\beta\Delta t}{\varepsilon}\right)$ . Therefore this successive penalty method is (at least) not worse than the Filbet-Jin scheme in stability. The resulting scheme is strongly AP, since  $f^n = M^n + \mathcal{O}(\varepsilon)$  for any  $n \geq 1$ , as  $\varepsilon \rightarrow 0$ . Besides, the positivity of  $f^{n+1}$  is preserved with the same technique and conditions introduced in section 2.1.2.

Compared with the Dimarco-Pareschi method, the exact solution for operator  $P$  is not needed. This scheme is applicable to a general  $P$ , as long as one can numerically solve the system implicitly. With  $Q(f)$  the Landau operator (3) and  $P(f)$  the Fokker-Planck operator (24), (44) gives a first order strongly AP scheme. Here  $\beta$  is given by (26).

**Remark 3.2** *On the computation cost.*

In the case of solving the Boltzmann equation by the BGK penalization, the three methods require the same amount of computation. The main cost is on the evaluation of the Boltzmann operator (2), which is solved by a fast spectral method proposed in [31, 15].

In the case of solving the Landau equation by the Fokker-Planck penalization, compared to the Filbet-Jin method, the successive penalty method requires one extra inversion of the Fokker-Planck operator at each  $x$  in every time step. However the cost of evaluating the Landau operator (3) is  $\mathcal{O}(N \log N)$  by the spectral method in [16]; while the cost of inverting the Fokker-Planck operator (24) is  $\mathcal{O}(N)$ , with a conjugate-gradient method (see [25] for details). In practice the computational cost does not increase significantly. As for the Dimarco-Pareschi method which requires the exact solution involving Fokker-Planck operator, the computation is much more costly than numerically solving the implicit system.

### 3.3 A second order successive penalty method

Both Filbet-Jin's and Dimarco-Pareschi's methods have second order extensions. Here we propose a second order successive method based on Filbet-Jin's method (28).

$$\begin{cases} \frac{f^* - f^n}{\Delta t/2} + v \cdot \nabla_x f^n = \frac{Q(f^n) - \beta^n P(f^n)}{\varepsilon} + \frac{\beta^n P(f^*)}{\varepsilon}, \\ \frac{f^{**} - f^n}{\Delta t} + v \cdot \nabla_x f^* = \frac{Q(f^*) - \beta^* P(f^*)}{\varepsilon} + \frac{\beta^*}{2\varepsilon} (P(f^n) + (1-\alpha)P(f^{**})), \\ \frac{f^{n+1} - f^{**}}{\Delta t} = \alpha \frac{\beta^* P(f^{n+1})}{\varepsilon}, \end{cases} \quad (49)$$

where  $P$  is the BGK operator (if  $Q$  is the Boltzmann operator) or the Fokker-Planck operator (if  $Q$  is the Landau operator). Note that a constant  $\alpha = \mathcal{O}(1)$  reduces the accuracy to first order. One has to take  $\alpha = \mathcal{O}(\Delta t)$ . For example, one can take

$$\alpha = \frac{\Delta t}{t_{\max}}. \quad (50)$$

This choice of splitting parameter is illustrated by the idea in [19]. Similar to the first order scheme, this choice of  $\alpha$  gives a stronger AP property as in (45). To have the strong AP property (7), one can take

$$\alpha = \alpha^n(x) = \min \left\{ \frac{\|f^n - M^n\|}{\Delta t}, \frac{\Delta t}{t_{\max}} \right\}. \quad (51)$$

Again the norm  $\|\cdot\|$  is taken over the velocity space.

**Remark 3.3** *One can also derive the corresponding split version of the second order method. A widely used technique is to apply the Strang splitting between the transport step and the collision step, with each step solved by a second order method separately. Here the transport part can be solved by a second order upwind method with slope limiters. The collision part can be solved by the same  $\mathcal{O}(\Delta t)$ -splitting successive penalization as in the nonsplit version. This indeed gives the second order accuracy for the Boltzmann equation 1 in the case  $\Delta t, \Delta x \ll \varepsilon$ . However, when  $\varepsilon \rightarrow 0$  while  $\Delta x$  and  $\Delta t$  fixed, one only obtains a first order method for the limit system (6) (see [21]). In other words, one cannot obtain a uniform second order method with Strang splitting.*

## 4 Numerical Tests

We always use the following settings, unless otherwise specified. The computation is performed on  $(x, v) \in [0, 1] \times [-v_{\max}, v_{\max}]^2$ , with  $v_{\max} = 8$ . We take  $N_x = 100$  grid points  $x$  direction and  $N_v = 32$  grid points in each  $v$  direction. We apply the van Leer type slope limiter [30] on the discretization of the transport parts, and take  $\Delta t = \frac{\Delta x}{2v_{\max}}$ , which guarantees the stability.

### 4.1 The AP property

We test the AP property of the first order Filbet-Jin method, Dimarco-Pareschi method and the successive penalty methods, in both split and nonsplit versions. The Boltzmann equation is solved with the BGK penalization. The AP properties for second order methods are similar.

The tests start with a non-equilibrium initial data,

$$f^0(x, v) = \frac{\rho^0(x)}{2\pi T^0(x)} \frac{1}{2} \left( e^{-\frac{|v-u^0(x)|^2}{2T^0(x)}} + e^{-\frac{|v+u^0(x)|^2}{2T^0(x)}} \right), \quad (52)$$

where

$$\begin{cases} \rho^0(x) = \frac{2 + \sin(2\pi x)}{3}, \\ u^0(x) = \begin{pmatrix} \cos(2\pi x) \\ 0 \end{pmatrix}, \\ T^0(x) = \frac{3 + \cos(2\pi x)}{4}. \end{cases} \quad (53)$$

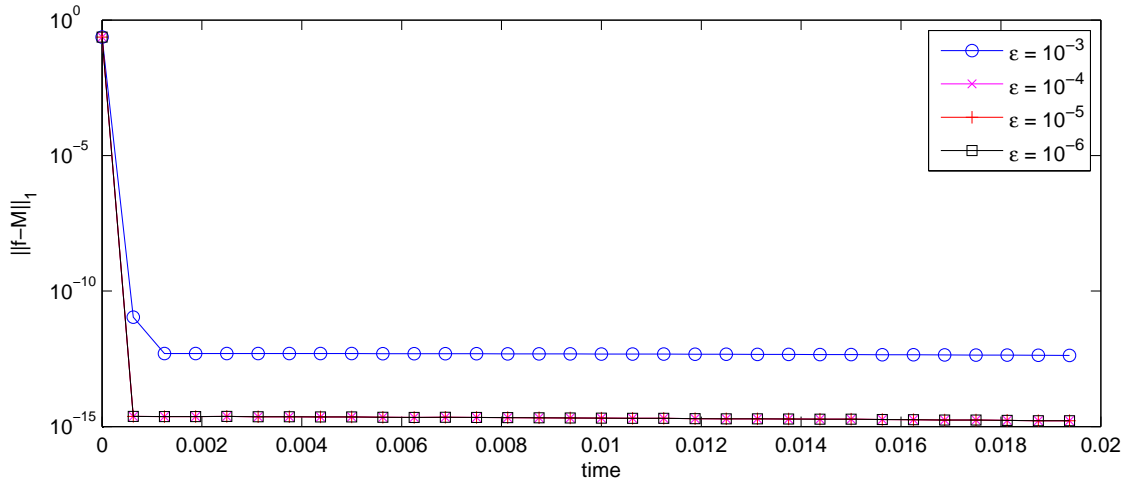
The periodic boundary condition is applied.

Figure 1 shows the time evolution of  $\|f - M\|_1$  (after the relaxation step) for the Dimarco-Pareschi method and the split version of the successive penalty method, for different  $\varepsilon$ . The solution from the Dimarco-Pareschi method has a much stronger compression effect on  $\|f - M\|_1$  as given in (10), while the solution of the successive penalty method has exactly the (strong) AP property we need.

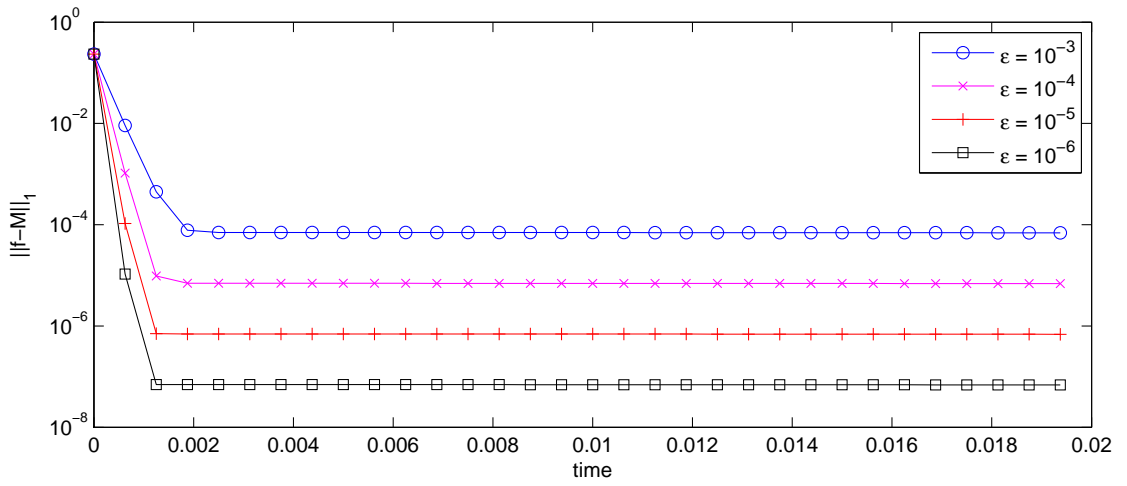
Figure 2 shows the time evolution of  $\|f - M\|_1$  for the Filbet-Jin method and the nonsplit version of successive method with  $\alpha$  given by (46) and  $\alpha = \frac{1}{2}$ , for different  $\varepsilon$ . The solution by the Filbet-Jin method shows the relaxed AP property, with the initial transient time, while the solution by the successive penalty method has the strong AP property we need. Note that with the constant  $\alpha = \frac{1}{2}$ , the method shows the over strong property

$$f^n - M^n = \mathcal{O}(\varepsilon^2).$$

We also give the AP results for the Landau equation. Figure 3 shows the time evolution of  $\|f - M\|_1$  for different methods, with different  $\varepsilon$ . As in the Boltzmann equation, the Filbet-Jin method gives a relaxed AP property, with an initial transient time. The split successive penalty method with  $\alpha = \frac{1}{2}$  and the nonsplit successive method with  $\alpha$  given by (46) show the strong AP property. The Dimarco-Pareschi method cannot be applied in this case.

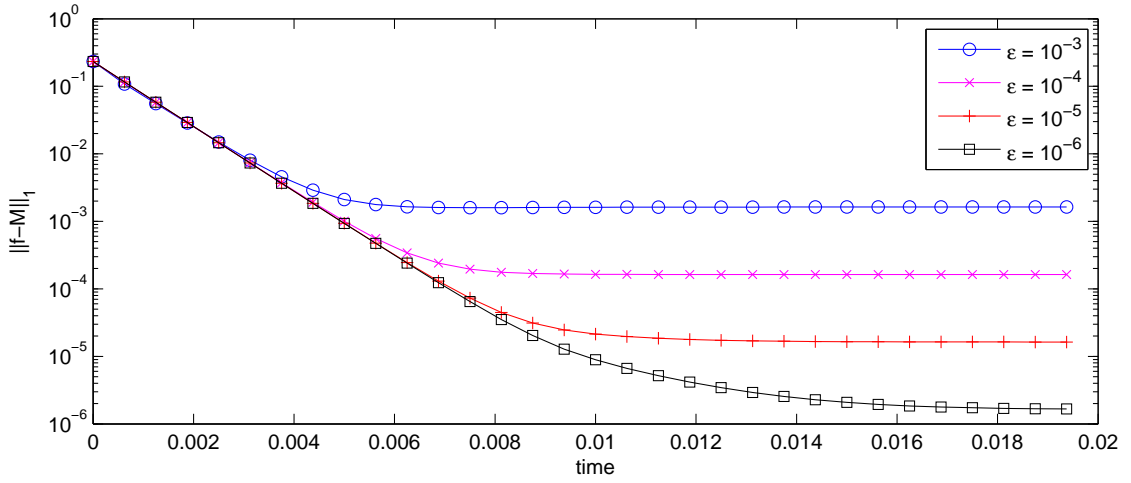


(a) The Dimarco-Pareschi method.

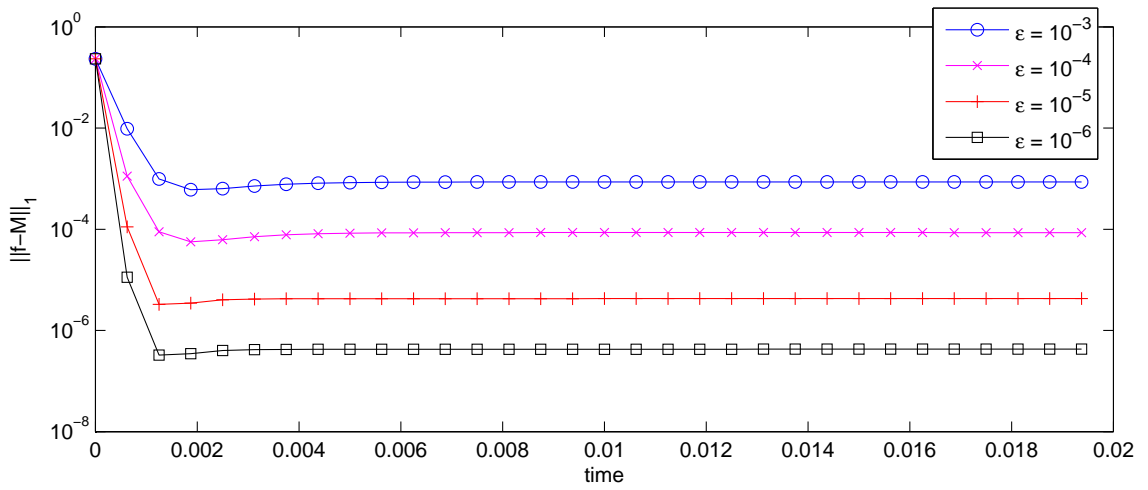


(b) The split successive penalty method.

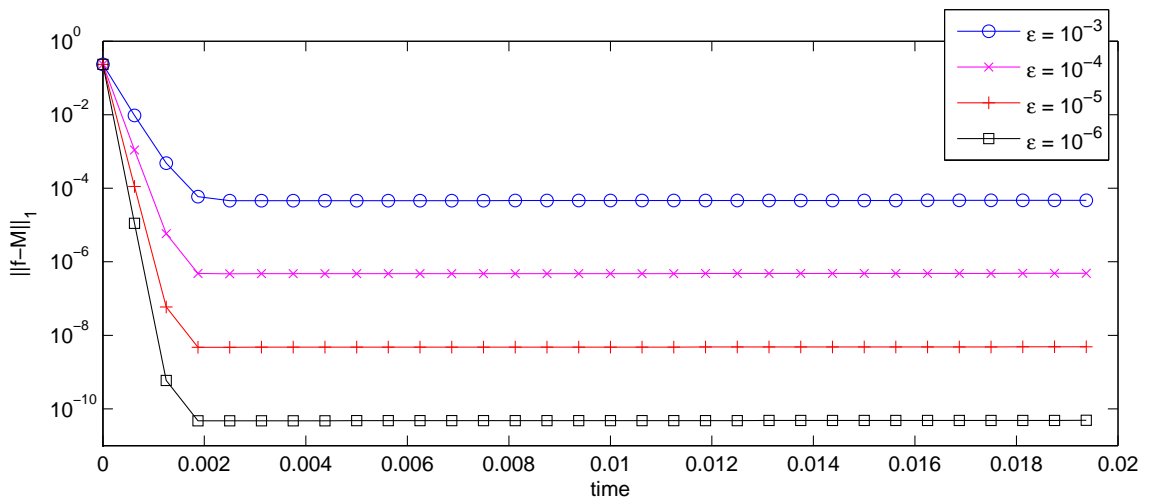
Figure 1: The time evolution of  $\|f - M\|_1$  for split methods for the Boltzmann equation.



(a) The Filbet-Jin method.

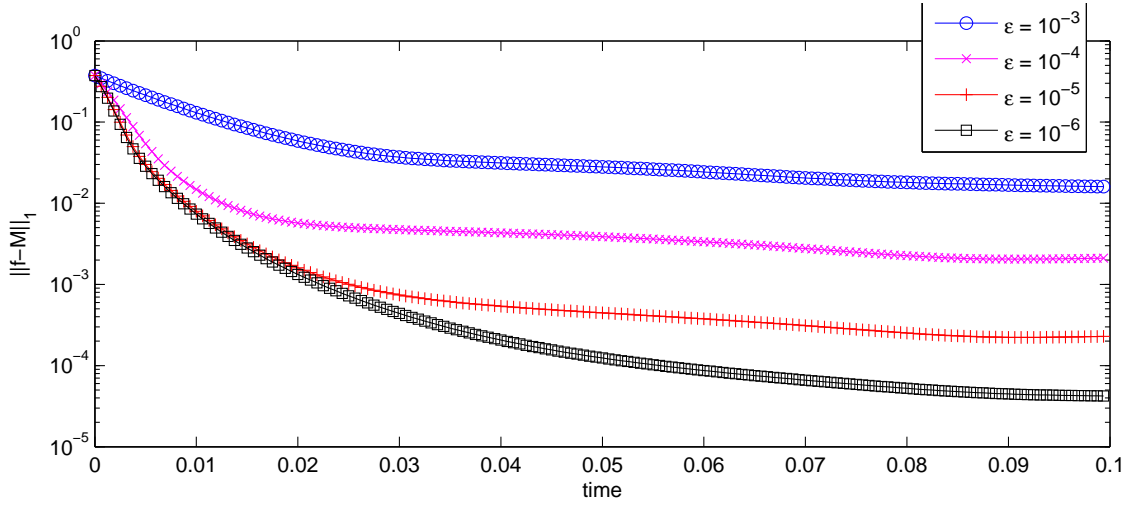


(b) The nonsplit successive penalty method with  $\alpha$  given by (46).

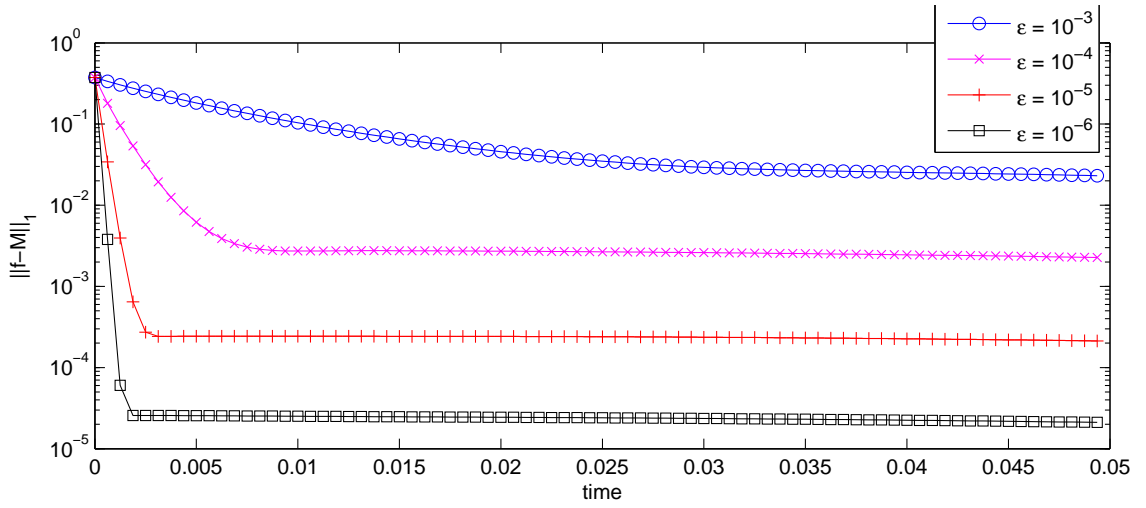


(c) The nonsplit successive penalty method with  $\alpha = 1/2$ .

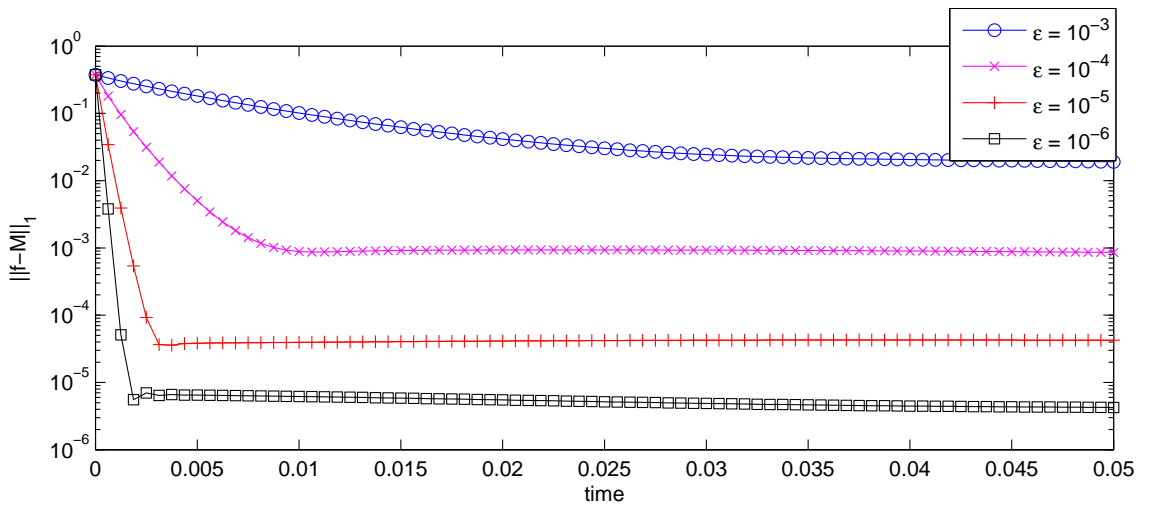
Figure 2: The time evolution of  $\|f - M\|_1$  for nonsplit methods for the Boltzmann equation.



(a) The Filbet-Jin method.



(b) The split successive penalty method with  $\alpha = 1/2$ .



(c) The nonsplit successive penalty method with  $\alpha$  by (46).

Figure 3: The time evolution of  $\|f - M\|_1$  for different methods with different  $\varepsilon$ , for the Landau equation.



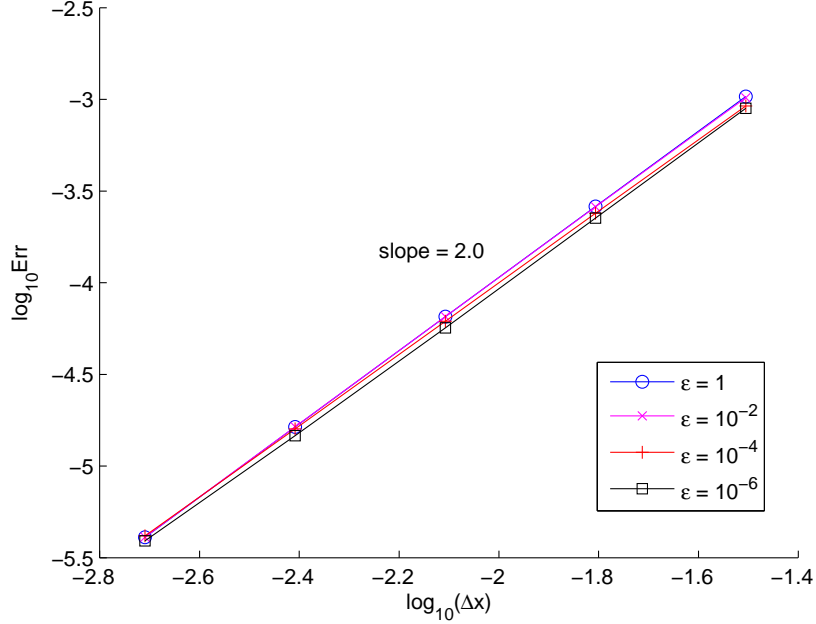


Figure 4: The test of convergence order for the successive penalty method with initial data (52)(53). This figure shows the  $l^1$  errors (54) with different  $\varepsilon$ .

## 4.2 Convergence order

Now we test the accuracy of the second order successive method (49) (51).

The non-equilibrium initial data (52)(53) are applied. We compute the solutions with grid points  $N_x = 32, 64, 128, 256, 512, 1024$  respectively. As mentioned before,  $N_v = 32$ . After time  $t_{\max} = 0.0625$  we check the following error,

$$e_{\Delta x}(f) = \max_{t \in (0, t_{\max})} \frac{\|f_{\Delta x}(t) - f_{2\Delta x}(t)\|_p}{\|f_{2\Delta x}(0)\|_p}. \quad (54)$$

This can be considered as an estimation of the relative error in  $l^p$  norm, where  $f_{\Delta x}$  are the numerical solutions computed from a grid of size  $\Delta x = \frac{1}{N_x}$ . The numerical scheme is said to be  $k$ -th order if  $e_{\Delta x} \leq C\Delta x^k$ , for  $\Delta x$  small enough.

Figure 4 gives the convergence order in  $l^1$  norm for the successive penalty method, with different  $\varepsilon$ . This shows that the scheme is second order in space (hence in time) **uniformly in  $\varepsilon$** , as expected.

## 4.3 The Riemann problem

Now we simulate the Sod shock tube problem, where the initial condition is  $f^I = M^I$  with

$$\begin{cases} (\rho, u_1, T) = (1, 0, 1), & \text{if } 0 \leq x < 0.5, \\ (\rho, u_1, T) = (1/8, 0, 1/4), & \text{if } 0.5 \leq x \leq 1. \end{cases} \quad (55)$$

The Neumann boundary condition in the  $x$ -direction is applied.

We apply the second order successive penalty method on this problem. We take  $N_x = 100$  and  $N_v = 32$ ,  $\Delta t = \frac{\Delta x}{2v_{\max}} \approx 6 \times 10^{-4}$ .

**Case I:**  $\varepsilon = 0.01$ .

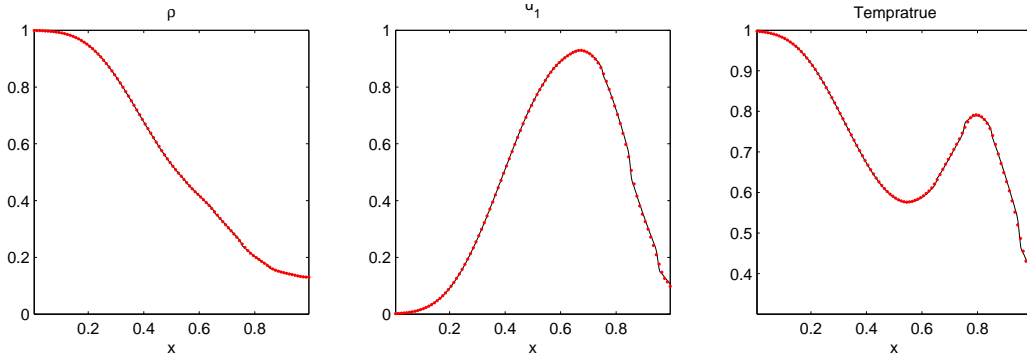


Figure 5: The comparison of density, velocity and temperature at  $t = 0.2$  between the resolved computation by the explicit second order Runger-Kutta scheme (solid line) and the under-resolved solutions by the second order successive penalty scheme (dots). Here  $\varepsilon = 0.01$ .

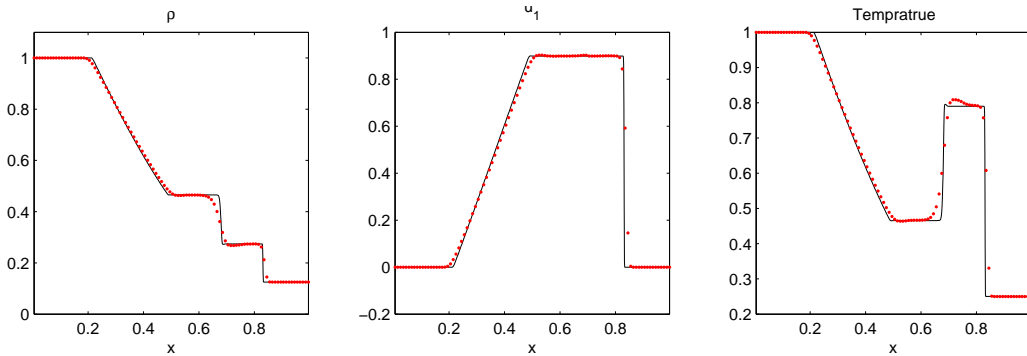


Figure 6: The comparison of density, velocity and temperature at  $t = 0.2$  between the solution of Euler system (solid line) and the under-resolved solutions by the second order successive penalty scheme (dots), with  $\varepsilon = 10^{-6}$ .

We compare this under-resolved solution to a fully resolved solution by the explicit second order Runger-Kutta scheme, where we take  $N_x = 1000$  and  $\Delta t = \frac{\Delta x}{2v_{\max}} \approx 6 \times 10^{-5}$ . We compute the macroscopic variables  $\rho$ ,  $u_1$  and  $T$ .

For such a value of  $\varepsilon$ , the problem is not stiff and this test is performed to compare the accuracy of our scheme with the classical (second order) RungeKutta method. The results are compared at  $t_{\max} = 0.2$  and shown in Figure 5. Therefore, in the kinetic regime our second order method gives the same accuracy as a second order fully explicit scheme without any additional computational effort.

**Case II:**  $\varepsilon = 10^{-6}$ .

Now the under-resolved solution is compared to the solution of the Euler system by a second order kinetic scheme, with  $N_x = 1000$  and  $\Delta t = 6 \times 10^{-5}$ . The macroscopic variables  $\rho$ ,  $u_1$  and  $T$  are compared at  $t_{\max} = 0.2$  and shown in Figure 6. The macroscopic quantities are well approximated although the mesh size and time steps are bigger than  $\varepsilon$ . The computational cost has been reduced significantly.

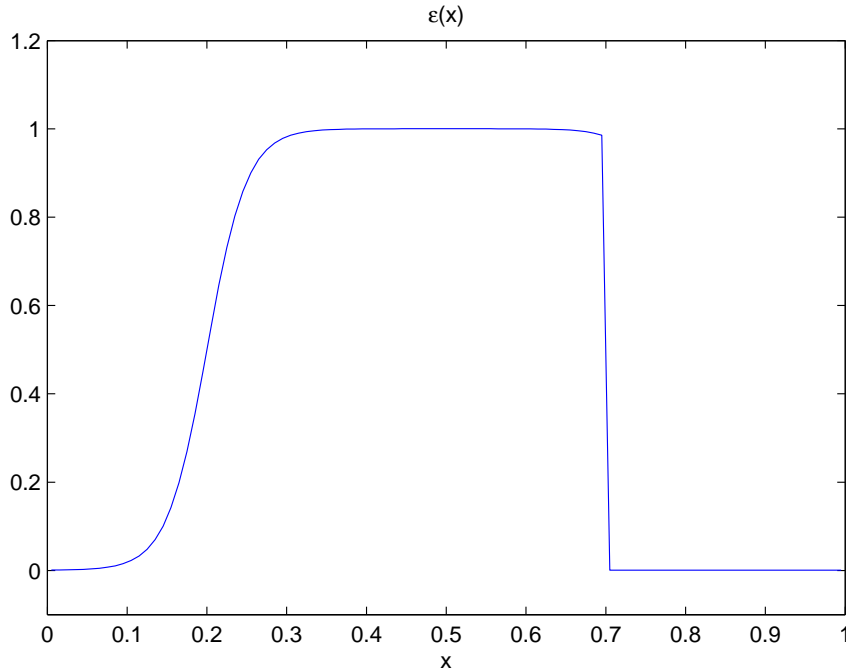


Figure 7: An  $x$ -dependent  $\varepsilon(x)$ .

#### 4.4 A mixing regime problem

Finally we apply the second order successive penalty method (49)(??) to the mixing regime problem ([14]). In this case the Knudsen number  $\varepsilon$  increases smoothly from  $\varepsilon_0$  to  $O(1)$ , then jumps back to  $\varepsilon_0$ ,

$$\varepsilon(x) = \begin{cases} \varepsilon_0 + \frac{1}{2} (\tanh(16 - 20x) + \tanh(-4 + 20x)), & x \leq 0.7 \\ \varepsilon_0, & x > 0.7 \end{cases}$$

with  $\varepsilon_0 = 0.0005$ . The picture of  $\varepsilon$  is shown in Figure 7. This problem involves mixed kinetic and fluid regimes.

To avoid the influence from the boundary, we take periodic boundary condition in  $x$ . The initial data are given by (52)(53).

In this test we compare the macroscopic variables obtained by our second order successive penalty scheme to the explicit Runger-Kutta scheme. For the explicit Runger-Kutta scheme, we take  $N_x = 1000$ ,  $\Delta t = \frac{\Delta x}{2v_{\max}} \approx 6 \times 10^{-5}$ . For our successive penalty scheme, we take  $N_x = 100$ ,  $\Delta t = \frac{\Delta x}{2v_{\max}} = 6 \times 10^{-4}$ . The results are compared up to  $t_{\max} = 0.75$  in Figure 8. Our scheme can capture the macroscopic behavior efficiently, with much larger mesh size and time steps.

## 5 Conclusions

In this paper we presented a successive penalty based asymptotic-preserving (AP) scheme for kinetic equations. This is an intermediate method between the Filbet-Jin method and the Dimarco-Pareschi method. It combines the advantages of both methods, with the same amount of computational cost. We presented a split version (42) and a nonsplit one (44), as well as their second order

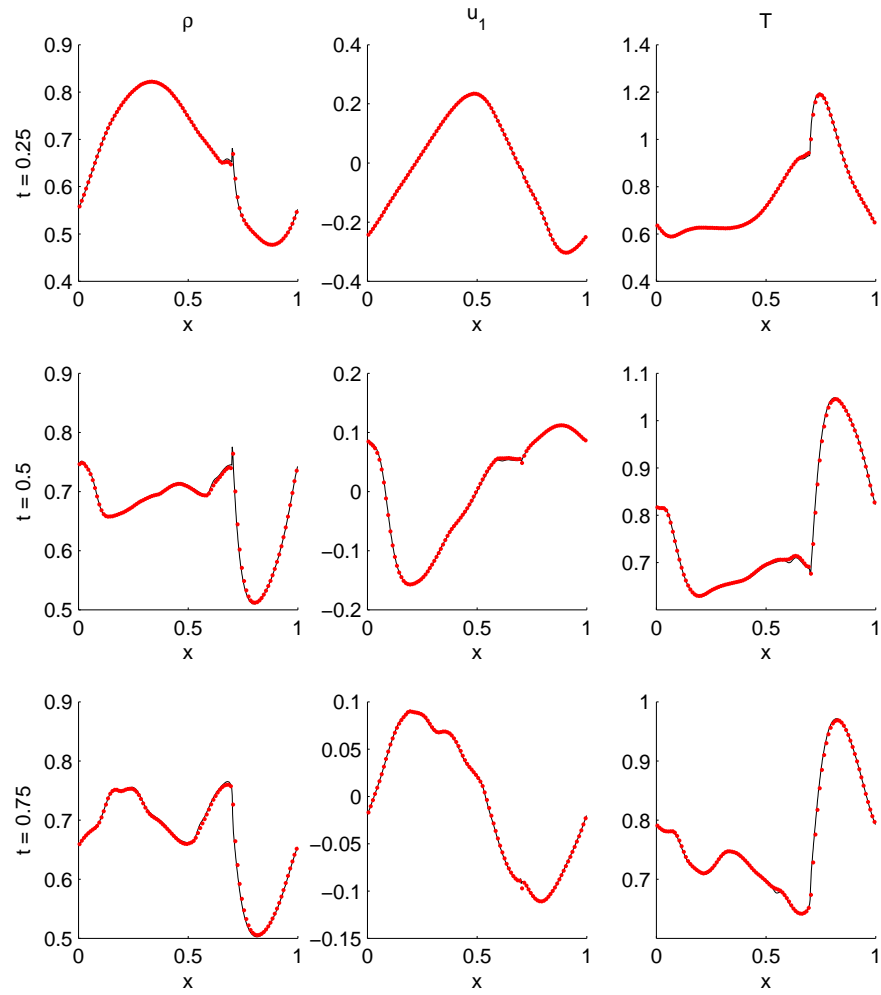


Figure 8: For mixing regime, the comparison between the resolved solutions (solid line) given by the explicit Runger-Kutta scheme and the solutions obtained by our successive penalty scheme (dots) with coarse grid and large time step.

extensions (??), (49). The new methods are strongly AP, positivity preserving, and applicable to very general collision operators, including the Boltzmann equation and the Landau equation.

## Acknowledgement

We thank a referee for a critical remark which helped to improve the paper.

## Appendix A The value of $\kappa$ in (22)

In this appendix we show it is not a problem to have a term  $M^{n+1}$ , which is close to 0 near the artificial boundary  $\{|v| = v_{\max}\} \subset V$ , in the denominator in the definition of  $\kappa$  (22). For this purpose, we only need to show the ratio of  $\frac{M^n}{M^{n+1}}$  does not “blow up” near the artificial boundary.

We give an estimation of  $\frac{M^n}{M^{n+1}}$ :

$$\begin{aligned} \frac{M^n}{M^{n+1}} &= \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \left\{ -\frac{(v-u^n)^2}{2T^n} + \frac{(v-u^{n+1})^2}{2T^{n+1}} \right\} \\ &= \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \left\{ -\frac{(v-u^n)^2}{2T^n} + \frac{(v-u^{n+1})^2}{2T^n} \right\} \exp \left\{ -\frac{(v-u^{n+1})^2}{2T^n} + \frac{(v-u^{n+1})^2}{2T^{n+1}} \right\} \\ &= \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \left\{ -\frac{(2v-u^n-u^{n+1})(u^{n+1}-u^n)}{2T^n} \right\} \exp \left\{ -\frac{(v-u^{n+1})^2(T^{n+1}-T^n)}{2T^n T^{n+1}} \right\}. \\ &= \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \left\{ -\frac{(2v-u^n-u^{n+1})\Delta t D_u}{2T^n} \right\} \exp \left\{ -\frac{(v-u^{n+1})^2 \Delta t D_T}{2T^n T^{n+1}} \right\}, \end{aligned}$$

with

$$D_u = \frac{u^{n+1} - u^n}{\Delta t}, \quad D_T = \frac{T^{n+1} - T^n}{\Delta t}.$$

Note that the CFL condition gives

$$\Delta t = \frac{\Delta x}{2v_{\max}},$$

therefore

$$\frac{M^n}{M^{n+1}} = \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \left\{ -\frac{D_u}{2T^n} \frac{(2v-u^n-u^{n+1})\Delta x}{2v_{\max}} \right\} \exp \left\{ -\frac{D_T}{2T^n T^{n+1}} \frac{(v-u^{n+1})^2}{2v_{\max}^2} v_{\max} \Delta x \right\}.$$

Noting  $\frac{|2v-u^n-u^{n+1}|}{2v_{\max}} \leq 2$  and  $\frac{(v-u^{n+1})^2}{2v_{\max}^2} \leq 2$ , the largest value is given by

$$\frac{M^n}{M^{n+1}} \leq \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} \exp \{C_1 \Delta x\} \exp \{C_2 v_{\max} \Delta x\},$$

where

$$C_1 = \frac{|D_u|}{T^n}, \quad C_2 = \frac{|D_T|}{T^n T^{n+1}}.$$

When  $\Delta x$  is small, one has

$$\frac{M^n}{M^{n+1}} = \frac{\rho^n T^{n+1}}{\rho^{n+1} T^n} (1 + O(\Delta x)).$$

In practice, we take  $v_{\max} = 8$  and  $\Delta x = \frac{1}{100}$ , the value of  $M^n/M^{n+1}$  is close to  $\frac{\rho^n T^{n+1}}{\rho^{n+1} T^n}$  in the whole domain.

In Figure 9, we give the numerical values of  $\beta$ ,  $\kappa$  and  $\beta(1 + \Delta t \kappa)$  in the first step of computation, with the initial data (52)(53). The value of  $\kappa$  is not very large, and  $\beta(1 + \Delta t \kappa)$  is very close to  $\beta$ . This gives an illustration of the typical values of these coefficients.

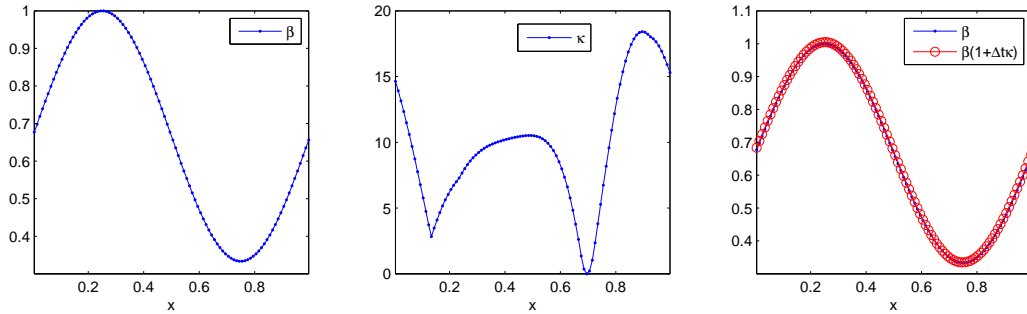


Figure 9: The values of  $\beta$ ,  $\kappa$  and  $\beta(1 + \Delta t\kappa)$  in the first step of computation, with the initial data (52)(53).

## References

- [1] O. B. A.A. ARSEN'EV, *On the connection between a solution of the Boltzmann equation and a solution of the Fokker-Planck-Landau equation*, Math. USSR Sbornik, 69 (1991), pp. 465–478.
- [2] C. BARDOS, F. GOLSE, AND D. LEVERMORE, *Fluid dynamic limits of kinetic equations. i. formal derivations*, Journal of Statistical Physics, 63 (1991), pp. 323–344.
- [3] M. BENNOUNE, M. LEMOU, AND L. MIEUSSENS, *Uniformly stable numerical schemes for the boltzmann equation preserving the compressible navier-stokes asymptotics*, Journal of Computational Physics, 227 (2008), pp. 3781 – 3803.
- [4] S. BOSCARINO, *Implicit-Explicit Runge-Kutta schemes for hyperbolic systems in the diffusion limit*, AIP Conference Proceedings, 1389 (2011), pp. 1315–1318.
- [5] F. BOUCHUT, F. GOLSE, AND M. PULVIRENTI, *Kinetic Equations and Asymptotic Theory*, Gauthiers-Villars, 2000.
- [6] C. CERCIGNANI, *The Boltzmann equation and its applications*, Springer-Verlag, 1988.
- [7] P. DEGOND AND B. LUCQUIN-DESREUX, *The Fokker-Planck asymptotics of the Boltzmann collision operator in the Coulomb case*, Mathematical Models and Methods in Applied Sciences (M3AS), 2 (1992), pp. 167–182.
- [8] J. DENG, *Asymptotic-preserving schemes for the semiconductor Boltzmann equation in the diffusive regime*, Numer. Math. Theor. Meth. Appl., 5 (2012), pp. 278–296.
- [9] L. DESVILLETES, *On asymptotics of the Boltzmann equation when the collisions become grazing*, Transport Theory and Statistical Physics, 21 (1992), pp. 259–276.
- [10] G. DIMARCO AND L. PARESCHI, *Exponential runge-kutta methods for stiff kinetic equations*, SIAM Journal on Numerical Analysis, 49 (2011), pp. 2057–2077.
- [11] E.GABETTA, L.PARESCHI, AND G.TOSCANI, *Wild's sums and numerical approximation of nonlinear kinetic equations*, Transport Theory and Statistical Physics, 25 (1996), pp. 515 – 531.
- [12] ———, *Relaxation schemes for nonlinear kinetic equations*, SIAM J. Numerical Analysis, 34 (1997), pp. 2168 – 2194.
- [13] F. FILBET, J. W. HU, AND S. JIN, *A numerical scheme for the quantum Boltzmann equation with stiff collision terms*, ESAIM-Math. Model. Numer. Anal., 46 (2012), pp. 443–463.
- [14] F. FILBET AND S. JIN, *A class of asymptotic-preserving schemes for kinetic equations and related problems with stiff sources*, J. Comp. Phys., 229 (2010), pp. 7625–7648.

- [15] F. FILBET, C. MOUHOT, AND L. PARESCHI, *Solving the Boltzmann equation in  $N \log N$* , SIAM J. Sci. Comput., 28 (2006), pp. 1029–1053.
- [16] F. FILBET AND L. PARESCHI, *A numerical method for the accurate solution of the Fokker-Planck-Landau equation in the nonhomogeneous case*, Journal of Computational Physics, 179 (2002), pp. 1 – 26.
- [17] F. FILBET AND A. RAMBAUD, *Analysis of an asymptotic preserving scheme for relaxation systems*, preprint.
- [18] T. GOUDON, *On Boltzmann equations and Fokker-Planck asymptotics: Influence of grazing collisions*, Journal of Statistical Physics, 89 (1997), pp. 751–776.
- [19] T. GOUDON, S. JIN, J. LIU, AND B. YAN, *Asymptotic-preserving schemes for kinetic-fluid modeling of disperse two-phase flows*, preprint.
- [20] J. HU, S. JIN, AND B. YAN, *A numerical scheme for the quantum Fokker-Planck-Landau equation efficient in the fluid regime*, Commun. Comput. Phys., 12 (2012), pp. 1541–1561.
- [21] S. JIN, *Runge-Kutta methods for hyperbolic conservation laws with stiff relaxation terms*, J. Comput. Phys., 122 (1995), pp. 51–67.
- [22] S. JIN, *Efficient asymptotic-preserving (AP) schemes for some multiscale kinetic equations*, SIAM J. Sci. Comput., 21 (1999), pp. 441–454.
- [23] ———, *Asymptotic preserving (AP) schemes for multiscale kinetic and hyperbolic equations: a review.*, Lecture Notes for Summer School on "Methods and Models of Kinetic Theory" (M&MKT), Porto Ercole (Grosseto, Italy), June 2010. Rivista di Matematica della Università di Parma, 3 (2012), pp. 177–216.
- [24] S. JIN AND Q. LI, *A BGK-penalization asymptotic-preserving scheme for the multispecies Boltzmann equation*, Numerical Methods for Partial Differential Equations, to appear.
- [25] S. JIN AND B. YAN, *A class of asymptotic-preserving schemes for the Fokker-Planck-Landau equation*, Journal of Computational Physics, 230 (2011), pp. 6420 – 6437.
- [26] L. LANDAU, *Die kinetische gleichung für den fall Coulombscher vechselwirkung*, Phys.Z. Sowjet, 154 (1963).
- [27] ———, *The transport equation in the case of the Coulomb interaction*, in Collected papers of L.D. Landau, D. ter Haar, ed., Pergamon press, Oxford, 1981, pp. 163–170.
- [28] M. LEMOU, *Relaxed micromacro schemes for kinetic equations*, C. R. Acad. Sci. Paris, Ser. I, 348 (2010), pp. 455–460.
- [29] M. LEMOU AND L. MIEUSSENS, *Implicit schemes for the Fokker-Planck-Landau equation*, SIAM Journal on Scientific Computing, 27 (2005), pp. 809–830.
- [30] R. LEVEQUE, *Numerical Methods for Conservation Laws*, Birkhauser-Verlag, Basel.
- [31] L. PARESCHI AND G. RUSSO, *Numerical solution of the Boltzmann equation I: Spectrally accurate approximation of the collision operator*, SIAM Journal on Numerical Analysis, 37 (2000), pp. pp. 1217–1245.
- [32] C. VILLANI, *A review of mathematical topics in collisional kinetic theory*, vol. 1 of Handbook of Mathematical Fluid Dynamics, North-Holland, 2002, pp. 71 – 74.