



**HAL**  
open science

# A FINITE VOLUME SCHEME FOR BOUNDARY-DRIVEN CONVECTION-DIFFUSION EQUATIONS WITH RELATIVE ENTROPY STRUCTURE

Francis Filbet, Maxime Herda

► **To cite this version:**

Francis Filbet, Maxime Herda. A FINITE VOLUME SCHEME FOR BOUNDARY-DRIVEN CONVECTION-DIFFUSION EQUATIONS WITH RELATIVE ENTROPY STRUCTURE. 2016. hal-01326029v2

**HAL Id: hal-01326029**

**<https://hal.science/hal-01326029v2>**

Preprint submitted on 5 Jul 2016 (v2), last revised 19 Apr 2017 (v4)

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

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

# A FINITE VOLUME SCHEME FOR BOUNDARY-DRIVEN CONVECTION-DIFFUSION EQUATIONS WITH RELATIVE ENTROPY STRUCTURE

FRANCIS FILBET AND MAXIME HERDA

ABSTRACT. We propose a finite volume scheme for a class of nonlinear parabolic equations endowed with non-homogeneous Dirichlet boundary conditions and which admit relative entropy functionals. For this kind of models including the porous media equations, Fokker-Planck equations for plasma physics or dumbbell models for polymer flows, it has been proved that the transient solution converges to a steady-state when time goes to infinity. The present scheme is built from the resolution of the stationary equation in order to preserve steady-states and natural Lyapunov functionals which provide a satisfying long-time behavior. After describing the numerical scheme, we present several numerical results which confirm the accuracy and underline the efficiency to preserve the large-time asymptotic.

KEYWORDS. Finite volume methods, relative entropy, non-homogeneous Dirichlet boundary conditions, polymers, magnetized plasma, porous media

MSC2010 SUBJECT CLASSIFICATIONS. 65M08, 65M12, 76S05, 76X05, 82D60

## CONTENTS

1. Introduction	1
1.1. General Setting	2
1.2. Physical models	4
1.3. Outline	5
2. Presentation of the numerical schemes	6
2.1. Discretization of the steady state equation	6
2.2. Discretization of evolution equation	7
2.3. Discrete relative $\phi$ -entropies	8
3. Analysis of the semi-discrete scheme	8
3.1. Relative entropy dissipation and stability	9
3.2. Proof of Theorem 3.1	12
3.3. Long-time behavior: proof of Theorem 3.2	12
4. Analysis of fully discrete schemes	13
4.1. Implicit Euler	13
4.2. Explicit Euler	14
5. Numerical simulations	16
5.1. Convergence and order of accuracy	16
5.2. Fokker-Planck with magnetic field	17
5.3. Polymer flow in a dilute solution	22
5.4. Porous medium equation	25
6. Comments and conclusion	26
References	27

## 1. INTRODUCTION

In this paper we propose to elaborate a finite volume scheme for nonlinear convection-diffusion equations with relative entropy structure set in a bounded domain and endowed with either non-homogeneous Dirichlet and/or null outward flux (generalized Neumann) boundary conditions. The

main objective of building such a scheme is to capture the correct long-time behavior when the solution converges to equilibrium.

**1.1. General Setting.** Let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$  with boundary  $\Gamma$ . Let us introduce an advection field  $\mathbf{E} : \Omega \rightarrow \mathbb{R}^d$  and  $\eta : \mathbb{R} \rightarrow \mathbb{R}$  a strictly increasing smooth function onto  $\mathbb{R}$  satisfying  $\eta(0) = 0$ . We consider the following nonlinear convection-diffusion equation with non-homogeneous Dirichlet boundary conditions

$$(1) \quad \begin{cases} \frac{\partial f}{\partial t} + \nabla \cdot (\mathbf{E} \eta(f) - \nabla \eta(f)) = 0 & \text{in } \mathbf{x} \in \Omega, \quad t \geq 0, \\ f = f^b & \text{on } \mathbf{x} \in \Gamma = \partial\Omega, \quad t \geq 0, \\ f(t=0) = f^{\text{in}} & \text{in } \mathbf{x} \in \Omega. \end{cases}$$

In [6], T. Bodineau, C. Villani, C. Mouhot and J. Lebowitz showed that this equation admits a large class of Lyapunov functionals, that we will denote, using their denomination, relative  $\phi$ -entropies. Each functional is generated by a convex function  $\phi$  and depends on a stationary state of (1). Therefore, we assume that there exists  $f^\infty$  which satisfies

$$(2) \quad \begin{cases} \nabla \cdot (\mathbf{E} \eta(f^\infty) - \nabla \eta(f^\infty)) = 0 & \text{in } \mathbf{x} \in \Omega, \\ f^\infty = f^b & \text{on } \mathbf{x} \in \Gamma. \end{cases}$$

Let us define the relative entropy corresponding to (1) and the associated dissipation.

**Definition 1.1** (Entropy generating functions). For any non-empty interval  $J$  of  $\mathbb{R}$  containing 1, we say that  $\phi \in \mathcal{C}^2(J, \mathbb{R}_+)$  is an *entropy generating function* or simply *entropy function* if it is strictly convex and satisfies  $\phi(1) = 0$  and  $\phi'(1) = 0$ .

Then from an entropy function  $\phi$ , the entropy functional is built as follows.

**Definition 1.2** (Relative  $\phi$ -entropy and dissipation). For any entropy generating function  $\phi$ , we denote by  $\mathcal{H}_\phi$  the so-called *relative  $\phi$ -entropy* defined by

$$\mathcal{H}_\phi(t) = \int_{\Omega} \int_{f^\infty(x)}^{f(t,x)} \phi' \left( \frac{\eta(s)}{\eta(f^\infty(x))} \right) ds \, d\mathbf{x},$$

and by  $\mathcal{D}_\phi$  the relative  $\phi$ -entropy dissipation defined by

$$\mathcal{D}_\phi(t) = \int_{\Omega} |\nabla h|^2 \phi''(h) \eta(f^\infty) \, d\mathbf{x},$$

where  $h$  is the ratio between the transient and stationary nonlinearities

$$(3) \quad h = \frac{\eta(f)}{\eta(f^\infty)}, \quad \text{if } \mathbf{x} \in \Omega \quad \text{and} \quad h = 1, \quad \text{if } \mathbf{x} \in \Gamma.$$

Typical examples of relative  $\phi$ -entropy are the so-called *physical relative entropies* and *p-entropies* (or *Tsallis relative entropies*) generated by, respectively,

$$(4) \quad \phi_1(x) = x \ln(x) - (x - 1), \quad \phi_p(x) = \frac{x^p - px}{p - 1} + 1 \quad \text{with } p \in (1, 2].$$

Let us note that for the linear problem, namely with  $\eta(x) = x$ , the relative  $\phi$ -entropy rewrites

$$\mathcal{H}_\phi(t) = \int_{\Omega} \phi \left( \frac{f}{f^\infty} \right) f^\infty \, d\mathbf{x}.$$

One readily sees that, since  $\eta$  and  $\phi'$  are increasing functions satisfying  $\eta(0) = 0$  and  $\phi'(1) = 0$ , the local relative  $\phi$ -entropy is a non-negative quantity. This yields  $H_\phi \geq 0$  and it cancels if and only if  $f$  and  $f^\infty$  coincide almost everywhere. The  $\phi$ -entropies are not, in general, distances between the solution and the steady state. However Csiszar-Kullback type inequalities [13, 22, 25] yield a control of the  $L^1$  distance between the solution and the equilibrium. Therefore if a relative  $\phi$ -entropy goes to zero when time goes to infinity, the solution converges to equilibrium in a strong sense.

The following result was proved in [6, Theorem 1.4] and yields the entropy-entropy dissipation principle for Equation (1), namely the decrease of all the relative  $\phi$ -entropies. The first step consists in reformulating the equation using the ratio (3). It writes

$$(5) \quad \frac{\partial f}{\partial t} + \nabla \cdot ([\mathbf{E}\eta(f^\infty) - \nabla\eta(f^\infty)] h - \eta(f^\infty)\nabla h) = 0.$$

**Proposition 1.3.** *Any  $L^\infty$  solution of (1) satisfies in the sense of distributions*

$$(6) \quad \frac{d\mathcal{H}_\phi}{dt} = -\mathcal{D}_\phi \leq 0,$$

for any entropy generating function  $\phi$ .

The formal computations leading to (6) motivate our choices in the elaboration of the discrete scheme. Therefore, we recall the proof yielding the entropy equality.

*Proof.* First, we integrate (5) against  $\phi'(h)$ , integrate by parts and use the boundary conditions and the fact that  $\phi'(1) = 0$  to get

$$\begin{aligned} \frac{d\mathcal{H}_\phi}{dt} &= \int_{\Omega} \phi'(h) \nabla \cdot [(\nabla\eta(f^\infty) - \mathbf{E}\eta(f^\infty)) h + \nabla h \eta(f^\infty)] \, d\mathbf{x} \\ &= - \int_{\Omega} (\nabla\eta(f^\infty) - \mathbf{E}\eta(f^\infty)) \cdot \nabla h \phi''(h) h \, d\mathbf{x} - \int_{\Omega} |\nabla h|^2 \phi''(h) \eta(f^\infty) \, d\mathbf{x}. \end{aligned}$$

Let  $\varphi$  be the only  $\mathcal{C}^1$  function satisfying,  $\varphi'(s) = \phi''(s) s$  and  $\varphi(1) = 0$ , given by  $\varphi(s) = s\phi'(s) - \phi(s)$ . Introducing it in the last expression yields

$$\begin{aligned} \frac{d\mathcal{H}_\phi}{dt} &= - \int_{\Omega} (\nabla\eta(f^\infty) - \mathbf{E}\eta(f^\infty)) \cdot \nabla\varphi(h) \, d\mathbf{x} - \mathcal{D}_\phi \\ &= -\mathcal{D}_\phi, \end{aligned}$$

where we integrated the first term by parts, used the stationary equation (2) and the boundary conditions.  $\square$

There are two important facts that justifies the use of (5) instead of (1) to derive the above entropy dissipation inequality. The rewriting transforms the advection field  $\mathbf{E}$  on  $\eta(f)$  into an incompressible field  $\nabla\eta(f^\infty) - \mathbf{E}\eta(f^\infty)$  on  $h$ . Therefore, the contribution of the convection vanishes when the time derivative of the relative entropy is computed and the convexity of  $\phi$  then suffices to provide dissipation. The underlying cancellations stems from the transformation of  $\nabla h \phi''(h) h$  into a gradient thanks to  $\varphi$  and on the fact that  $f^\infty$  solves (2). The second reason is that considering the equation on  $h$  instead of  $f$  changes non-homogeneous Dirichlet boundary conditions into homogeneous ones on  $h - 1$ . This and the properties of  $\phi$  enables the cancellations of boundary terms. In other words, the  $\phi$  relative entropies are the correct functionals and (5) the right form of the equation to capture the boundary-driven dynamic. The purpose of this work is the preservation by a discrete scheme of the *whole class* of relative entropy dissipation inequalities of the continuous model. This will be done adapting the above strategy to a finite volume discretization of the equation.

Let us emphasize that  $\mathbf{E}$  is a general field and need not to be either incompressible nor irrotational as for parabolic equations with a gradient flow structure [24]. Indeed, assuming some regularity on the advection field, one can apply the Hodge decomposition to get the existence of a potential  $\varphi : \Omega \rightarrow \mathbb{R}$  and  $\mathbf{F} : \Omega \rightarrow \mathbb{R}^d$  such that

$$(7) \quad \mathbf{E} = \nabla\varphi + \mathbf{F}, \quad \nabla \cdot \mathbf{F} = 0$$

When  $\mathbf{F} = \mathbf{0}$ , there are many examples in the literature [19, 11, 5, 9, 8, 12] of finite volume schemes preserving entropy dissipation properties. C. Chainais-Hillairet and F. Filbet studied in [11] a finite volume discretization for nonlinear drift-diffusion system and proved that the numerical solution converges to a steady-state when time goes to infinity. In [8], M. Burger, J. A. Carrillo and M. T. Wolfram proposed a mixed finite element method for nonlinear diffusion equations and proved convergence towards the steady-state in case of a nonlinear Fokker-Planck equation with uniformly convex potential. All these schemes exploit the gradient flow structure of the equation, which gives a natural entropy.

In the non-symmetric case  $\mathbf{F} \neq \mathbf{0}$ , the gradient structure cannot be exploited anymore, but as Proposition 1.3 show, there is still a relative entropy structure, which may be investigated to prove convergence to a steady state. The relative entropy properties of Fokker-Planck type equations in the whole space are exhaustively studied in the famous paper [3] of A. Arnold, P. Markowich, G. Toscani, A. Unterreiter and specific properties of the non-symmetric equations have been investigated in [2, 1].

In bounded domains, entropy properties are often used in the context of no-flux boundary conditions or in the whole space, but few results concern Dirichlet boundary conditions. In [5], M. Bessemoulin-Chatard proposed an extension of the Scharfetter-Gummel for finite volume scheme for convection-diffusion equations with nonlinear diffusion and non-homogeneous and unsteady Dirichlet boundary conditions. While in the latter work the author presents a scheme with a satisfying long-time behavior for a larger class of models than those of the present paper, our strategy and objectives differ. Here we aim at preserving a whole class of relative entropies and build our scheme for the transient problem from a discretization of the stationary equation.

Let us precise that we can generalize our approach to the more general boundary conditions

$$\begin{cases} f = f^b \text{ on } \Gamma_D, \\ [\mathbf{E}\eta(f) - \nabla\eta(f)] \cdot \mathbf{n}(\mathbf{x}) = 0 \text{ on } \Gamma_N, \end{cases}$$

with  $\Gamma = \Gamma_D \cup \Gamma_N$ . Our results hold in this setting with minor modifications but to avoid unnecessary technicalities in the notation and in the analysis we consider non-homogeneous Dirichlet conditions on the whole boundary. However, numerical results will be shown in both cases.

**1.2. Physical models.** Before describing our numerical scheme, let us present some physical models described by equation (1) for which the large-time asymptotic has been studied using entropy/entropy-dissipation arguments. Some of these models are the homogeneous part of kinetic Fokker-Planck-type equations and this work constitutes a first step towards treating full kinetic models. In future work, we aim at adapting the strategy developed here to ensure the property of convergence to local equilibrium for the solutions of these equations.

**1.2.1. The Fokker-Planck equation with magnetic field.** A classical model of plasma physics describing the dynamic of charged particles evolving in an external electromagnetic field  $(-\nabla_{\mathbf{x}}\phi(\mathbf{x}), \mathbf{B}(\mathbf{x}))$  is given by the Vlasov-Fokker-Planck equation reading

$$(8) \quad \frac{\partial F}{\partial t} + \mathbf{v} \cdot \nabla_{\mathbf{x}} F - \nabla_{\mathbf{x}} \phi \cdot \nabla_{\mathbf{v}} F + (\mathbf{v} \wedge \mathbf{B}) \cdot \nabla_{\mathbf{v}} F = \nabla_{\mathbf{v}} \cdot (\mathbf{v} F + \nabla_{\mathbf{v}} F).$$

For more details on the model we refer to [7, 20]. In [7], Bouchut and Dolbeault proved that the solution of (8) in the whole phase space and without magnetic field converges to a global equilibrium. Their proof mainly relies on the decrease of the free energy functional, which corresponds to the physical relative entropy introduced in (4). The external magnetic field does not alter the relative entropy inequality. We refer to [20] for the corresponding computations. Here, we consider the phenomena happening in the velocity space which results in a Fokker-Planck equation with magnetic field, namely equation (1) with  $\eta(s) = s$  and where the advection field is given by

$$(9) \quad \mathbf{E}(\mathbf{v}) = -\mathbf{v} + \mathbf{v} \wedge \mathbf{B},$$

with constant magnetic field  $\mathbf{B}$ . In applications, the velocity variable  $\mathbf{v}$  usually lives in  $\mathbb{R}^3$ . However when performing numerical simulations, one needs to restrict the velocity domain to a bounded set  $\Omega$ . On the edge of this restricted domain  $\Omega$ , we shall consider the following non-homogeneous Dirichlet boundary conditions

$$(10) \quad f(t, \mathbf{v}) = f^\infty(\mathbf{v}) \quad \forall \mathbf{v} \in \partial\Omega,$$

where  $f^\infty$  is the local Maxwellian associated with (8) which writes

$$(11) \quad f^\infty(\mathbf{v}) = \frac{1}{(2\pi)^{3/2}} e^{-\frac{|\mathbf{v}|^2}{2}} \quad \forall \mathbf{v} \in \Omega,$$

and is a stationary state of (1)-(9)-(10). With the boundary conditions (10), one recovers the same stationary state as in the whole space while working in a bounded domain. As for the more complicated kinetic model (8), free energy (relative  $\phi$ -entropy) decrease holds.

Our approach is particularly promising for this kind of problem when the solution develops some micro-instabilities around a steady state. In this situation, it is important that numerical artifacts do not generate some spurious oscillations.

**1.2.2. The dumbbell model for the density of polymers in a dilute solution.** The following kinetic equation describes the evolution of the density  $F(t, \mathbf{x}, \mathbf{k})$  of polymers at time  $t$  and position  $\mathbf{x}$  diluted in a fluid flow of velocity  $\mathbf{u}(\mathbf{x})$  from a mesoscopic point of view

$$(12) \quad \frac{\partial F}{\partial t} + \mathbf{u} \cdot \nabla_{\mathbf{x}} F = -\nabla_{\mathbf{k}} \cdot \left[ \left( \nabla_{\mathbf{x}} \mathbf{u} \mathbf{k} - \frac{1}{2} \nabla_{\mathbf{k}} \Pi(\mathbf{k}) \right) F - \frac{1}{2} \nabla_{\mathbf{k}} F \right].$$

The polymers are pictured as two beads linked by a spring and the variable  $\mathbf{k}$  stands for the vector indicating the length and orientation of the molecules. The potential  $\Pi$  is given by  $\Pi(\mathbf{k}) = |\mathbf{k}|^2/2$  in the case of Hookean dumbbells or by  $\Pi(\mathbf{k}) = -\ln(1-|\mathbf{k}|^2)/2$  in the case of Finite Extensible Nonlinear Elastic dumbbells. In the complete model, the velocity of the fluid  $\mathbf{u}$  follows an incompressible Navier-Stokes equation featuring an additional force term modeling for the contribution of the polymers on the dynamic of the fluid which results in a nonlinear kinetic-fluid coupling. Here, we consider the simpler case where  $u(\mathbf{x})$  is a given incompressible field. For more details on the modeling, we refer to [21] and references therein. We also refer to the paper [23] of Masmoudi that treats the well-posedness and provides additional information on the model.

Once again we aim at approximating numerically the "velocity" part of the kinetic equation (12) which rewrites as (1) with  $\eta(s) = s$  and where the advection field is given by

$$(13) \quad \mathbf{E}(\mathbf{k}) = \mathbf{A} \mathbf{k} - \frac{1}{2} \nabla_{\mathbf{k}} \Pi(\mathbf{k}),$$

with a constant matrix  $\mathbf{A}$  satisfying  $\text{tr}(\mathbf{A}) = 0$  and to be seen as the gradient of an incompressible velocity field at some space location. Natural boundary conditions for this model are given by null outward flux. In [21], the long-time behavior of (12) and of the latter reduced model are investigated using relative  $\phi$ -entropies where  $\phi$  may typically be given by (4).

**1.2.3. A nonlinear model, the porous medium equation.** The porous medium equation is a nonlinear PDE writing

$$(14) \quad \frac{\partial f}{\partial t} = \Delta f^m,$$

with  $m > 1$ . It can model many physical applications and generally describes processes involving fluid flow, heat transfer or diffusion. The typical example is the description of the flow of an isentropic gas through a porous medium. There is a huge literature on this equation and we refer to the book of Vázquez [26] for the detailed mathematical theory.

Here, equation (14) is set in a bounded domain  $\Omega$  with non-homogeneous Dirichlet boundary conditions

$$f(t, \mathbf{x}) = f^b(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in \partial\Omega,$$

such that it might be recast like (1) with a null advection field and with  $\eta(s) = s^m$ . Using their relative  $\phi$ -entropy method, Bodineau, Mouhot, Villani and Lebowitz show exponential convergence to equilibrium for this nonlinear equation.

**1.3. Outline.** The plan of the paper is as follows. In Section 2, we present the semi-discrete finite volume scheme and the discrete version of the relative  $\phi$ -entropies. Then we prove in Section 3 the main properties of our scheme, namely the relative  $\phi$ -entropy dissipation as well as well-posedness, stability and long-time behavior of solutions. In Section 4, we adapt these results to the implicit and explicit Euler time discretization of the scheme. We end up providing numerical proofs of the properties of our schemes as well as a convergence analysis on several test cases consisting in three-dimensional versions of models presented above, with their various boundary conditions.

## 2. PRESENTATION OF THE NUMERICAL SCHEMES

In this section, we introduce our scheme in the semi-discrete form. We start with some notation associated to our finite volume approximation of Equation (5).

A discretization of  $\Omega$ , is defined by the triplet  $\mathcal{D} = (\mathcal{T}, \mathcal{E}, \mathcal{P})$ . The mesh  $\mathcal{T}$  is a finite family of nonempty connected open disjoint subsets of  $\Omega$  called the control volumes  $K \in \mathcal{T}$ . The closure of the union of all the control volumes shall be equal to  $\bar{\Omega}$ . The set  $\mathcal{E}$  is a finite family of nonempty open disjoint subsets of  $\bar{\Omega}$  called the edges  $\sigma \in \mathcal{E}$ . Each edge is a subset of an affine hyperplane in  $\mathbb{R}^{d-1}$  with positive measure. It is also assumed that for any control volume  $K \in \mathcal{T}$  there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that the closure of the union of all the edges in  $\mathcal{E}_K$  is equal to  $\partial K = \bar{K} \setminus K$ . We also define several subsets of  $\mathcal{E}$ . The family of interior edges  $\mathcal{E}_{\text{int}}$  is given by  $\{\sigma \in \mathcal{E}, \sigma \not\subseteq \Gamma\}$  and the family of exterior edges by  $\mathcal{E}_{\text{ext}} = \mathcal{E} \setminus \mathcal{E}_{\text{int}}$ . Similarly, for any control volume  $K \in \mathcal{T}$ , we define  $\mathcal{E}_{\text{int},K} = \mathcal{E}_{\text{int}} \cap \mathcal{E}_K$  and  $\mathcal{E}_{\text{ext},K} = \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ . Moreover we assume that for any edge  $\sigma$ , the number of control volumes sharing the edge  $\sigma$ , namely the cardinality of  $\{K \in \mathcal{T}, \sigma \in \mathcal{E}_K\}$  is exactly 2 for interior edges and 1 for exterior edges. With this assumptions, every interior edge is shared by two control volumes, say  $K$  and  $L$ , so that we may use the notation  $\sigma = K|L$  whenever  $\sigma \in \mathcal{E}_{\text{int}}$ . The set  $\mathcal{P} = \{\mathbf{x}_K\}_{K \in \mathcal{T}}$  is a finite family of points satisfying that for any control volume  $K \in \mathcal{T}$ ,  $\mathbf{x}_K \in K$ . The Dirichlet condition on the boundary is given by  $f^b \in L^\infty(\Gamma)$  which is assumed to be positive. The approximate solution  $f(t)$  at time  $t \in [0, T)$  is an element of the set

$$X_{fb} = \left\{ f \in \mathbb{R}^{\mathcal{T}} \times \mathbb{R}^{\mathcal{E}_{\text{ext}}} : f_\sigma = \frac{1}{m(\sigma)} \int_\sigma f^b dm, \forall \sigma \in \mathcal{E}_{\text{ext}} \right\}.$$

For any function  $\psi : \mathbb{R} \rightarrow \mathbb{R}$ , and  $f \in X_{fb}$  we shall define the component-wise composition with the intuitive notation  $\psi(f) = ((\psi(f_K))_{K \in \mathcal{T}}, (\psi(f_\sigma))_{\sigma \in \mathcal{E}_{\text{ext}}})$ . Finally, to ease the notation we sometimes denote by  $f$  the piecewise constant function satisfying  $f(x) = f_K$  almost everywhere for each  $K \in \mathcal{T}$ .

In order to get a numerical approximation of the solution of (1), for which the numerical solution converges to a consistent steady state of (2) and satisfies the discrete equivalent of the entropy-entropy dissipation equality (6), we proceed in two steps.

- We first solve a discrete steady state problem consistent with (2).
- We use the steady state to define a numerical flux such that the numerical solution satisfies discrete equivalents of the  $\phi$ -entropy inequalities (6). In particular we want the numerical solution to converge to the discrete steady state when time goes to infinity.

**2.1. Discretization of the steady state equation.** In this subsection, we look for an approximation  $f^\infty \in X_{fb}$  of the continuous stationary state (2). This means that we construct a discrete flux which approximates consistently the flux  $(\mathbf{E}\eta(f^\infty) - \nabla\eta(f^\infty)) \cdot \mathbf{n}$  of the stationary equation (2) and such that its discrete divergence cancels, namely

$$(15) \quad \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^\infty = 0.$$

Our method do allow to approach the fluxes from the analytical steady states as well as finding an approximate steady state from a finite volume discretization. For some models, the global equilibrium  $f^\infty$  may be known analytically. Therefore, we construct a discrete approximation in  $X_{fb}$  by a standard projection

$$f_K^\infty = \frac{1}{m(K)} \int_K f^\infty d\mathbf{x},$$

and the numerical flux  $F_{K,\sigma}^\infty$  is given by

$$F_{K,\sigma}^\infty = \int_\sigma (\mathbf{E}\eta(f^\infty) - \nabla\eta(f^\infty)) \cdot \mathbf{n}_{K,\sigma} dm,$$

where the integrals may be computed exactly or approximated with a quadrature formula. In any case, equation (15) shall be satisfied as well as the interior continuity condition,

$$F_{K,\sigma}^\infty = -F_{L,\sigma}^\infty, \quad \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}.$$

Let us precise that this last condition is required and satisfied for all the discrete flux we define in the following.

On the other hand, when the steady state is not explicitly known we may apply a finite volume scheme to compute a numerical approximation. Once again, *our method do not impose any particular scheme for solving the stationary equation* but for the sake of completeness we propose the following that we actually use in some of our numerical simulations. Keeping in mind that most of the models we consider arise as space-homogeneous part of kinetic equations, we need to deal with potentially small (but non-zero) Dirichlet boundary conditions as well as confining potentials. this results in Maxwellian-like equilibria where  $f^\infty$  is maximal somewhere inside the domain and quickly decays towards small values at the boundary. Therefore finding a scheme that can provide a good approximation with this few information on the boundary is a hard task in itself. Here is our proposition to deal with this issue.

We replace the unknown  $f^\infty$  by  $h^\infty = \eta(f^\infty)/\exp(\phi)$  where  $\phi$  is the potential of the advection field  $\mathbf{E}$  in its Hodge decomposition (7). At the continuous level the stationary equation (2) then rewrites

$$\begin{cases} \nabla \cdot [e^\phi(\mathbf{F} h^\infty - \nabla h^\infty)] = 0 & \text{in } \mathbf{x} \in \Omega, \\ h^\infty = \eta(f^b) \exp(-\phi) & \text{on } \mathbf{x} \in \Gamma. \end{cases}$$

With the above formulation of Equation (2), and provided that  $f^b$  is of the same order as  $\exp(\phi)$  on the boundary, we expect the solution  $h^\infty$  to be close to 1 inside the domain and on the boundary. The corresponding finite volume scheme is given by (15), where the fluxes  $F_{K,\sigma}^\infty$  are discretized with upwind discretization for the convective part and with a two-points centered gradient for the diffusion.

**2.2. Discretization of evolution equation.** Now we treat the time evolution problem and use the numerical flux constructed from the steady state problem (2) to build a numerical approximation of the time evolution equation. The scheme reads in semi-discrete form

$$(16) \quad \forall K \in \mathcal{T}, \quad \begin{cases} m(K) \frac{df_K}{dt}(t) + \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma} = 0, & \forall t \in [0, T), \\ f_K(0) = \frac{1}{m(K)} \int_K f^{\text{in}}(\mathbf{x}) \, d\mathbf{x}. \end{cases}$$

We first define  $h \in X_1$  as  $h = \eta(f)/\eta(f^\infty)$  with component-wise division and for each  $K \in \mathcal{T}$  and  $\sigma \in \partial K$ , the flux is given by

$$(17) \quad F_{K,\sigma} = F_{K,\sigma}^{\text{conv}} + F_{K,\sigma}^{\text{diss}},$$

where we call  $F_{K,\sigma}^{\text{conv}}$  the convective flux and  $F_{K,\sigma}^{\text{diss}}$  the dissipative flux corresponding to the diffusive term. The convective flux approximates the advection part of the continuous flux in (5), namely  $[\mathbf{E} \eta(f^\infty) - \nabla \eta(f^\infty)] h \cdot \mathbf{n}$ . The corresponding discrete velocity field is given by

$$U_{K,\sigma} = \frac{1}{m(\sigma)} F_{K,\sigma}^\infty,$$

and we define the associated monotone convective flux by

$$(18) \quad F_{K,\sigma}^{\text{conv}} = \begin{cases} m(\sigma) [U_{K,\sigma}^+ g(h_K, h_L) - U_{K,\sigma}^- g(h_L, h_K)] & \text{if } \sigma = K|L, \\ m(\sigma) [U_{K,\sigma}^+ g(h_K, h_\sigma) - U_{K,\sigma}^- g(h_\sigma, h_K)] & \text{otherwise,} \end{cases}$$

where for any real number  $u$ , we define the positive and negative part of  $u$  to be respectively  $u^+ = \max(u, 0)$  and  $u^- = \max(-u, 0)$ . The function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  is locally Lipschitz-continuous, non-decreasing in the first variable, non-increasing in the second variable and satisfies  $g(s, s) = s$  for consistency. In the numerical simulations, we use a classical upwind flux which corresponds to  $g(s, t) = s$ . The dissipative flux  $F_{K,\sigma}^{\text{diss}}$  is an approximation of the diffusion part of the continuous flux



$-\eta(f^\infty)\nabla h \cdot \mathbf{n}$  in (5) and is built on a standard centered approximation of the derivative along the outward normal vector of each edge, namely

$$(19) \quad F_{K,\sigma}^{\text{diss}} = -\tau_\sigma \eta(f_\sigma^\infty) D_{K,\sigma} h,$$

where  $f_\sigma^\infty$  is a consistent approximation of the stationary state on the edge  $\sigma$  to be chosen. In the following we just suppose that it is given by, say,  $f_\sigma^\infty = (f_K^\infty + f_L^\infty)/2$ . The quantity  $\tau_\sigma$  is the transmissibility of the edge  $\sigma$ , given by

$$\tau_\sigma = \frac{m(\sigma)}{d_\sigma},$$

where

$$d_\sigma = \begin{cases} d(\mathbf{x}_K, \sigma) + d(\mathbf{x}_L, \sigma) & \text{if } \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \\ d(\mathbf{x}_K, \sigma) & \text{if } \sigma \in \mathcal{E}_{\text{ext},K}, \end{cases}$$

with  $d(\cdot, \cdot)$  the euclidean distance  $\mathbb{R}^d$ . The difference operator  $D_{K,\sigma}$  is defined for any  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$  by  $D_{K,\sigma} : \mathbb{R}^{\mathcal{T}} \times \mathbb{R}^{\mathcal{E}_{\text{ext}}K} \rightarrow \mathbb{R}$ ,

$$(20) \quad D_{K,\sigma} u = \begin{cases} u_L - u_K & \text{if } \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \\ u_\sigma - u_K & \text{if } \sigma \in \mathcal{E}_{\text{ext},K}. \end{cases}$$

For consistency of the discrete gradients, we require an orthogonality condition for the mesh, namely

$$\forall \mathbf{x}, \mathbf{y} \in \sigma = K|L, \quad (\mathbf{x} - \mathbf{y}) \cdot (\mathbf{x}_K - \mathbf{x}_L) = 0.$$

**2.3. Discrete relative  $\phi$ -entropies.** For  $f \in X_{fb}$ , the semi-discrete equivalent of the relative  $\phi$ -entropy in (1.2) is given by

$$(21) \quad H_\phi(f) = \sum_{K \in \mathcal{T}} m(K) e_{\phi,K}(f),$$

where  $e_\phi = (e_{\phi,K})_{K \in \mathcal{T}}$  is the local relative  $\phi$ -entropy writing

$$e_{\phi,K}(f) = \int_{f_K^\infty}^{f_K} \phi' \left( \frac{\eta(s)}{\eta(f_K^\infty)} \right) ds, \quad \forall K \in \mathcal{T}.$$

We also introduce the contributions of the convective and diffusive part of the equation to the relative entropy variation

$$(22) \quad C_\phi(f) = \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \phi'(h_K) F_{K,\sigma}^{\text{conv}}, \quad D_\phi(f) = \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \phi'(h_K) F_{K,\sigma}^{\text{diss}}.$$

In Proposition 3.3 we will show that  $D_\phi$  is consistent with its continuous analogue  $\mathcal{D}_\phi$ . In the continuous setting, the contribution of the advection vanishes in the variation of the relative  $\phi$ -entropy. We will show in Proposition 3.3 that it is not the case in the discrete setting. However the monotonicity properties of the convective flux make this term an additional numerical dissipation to the relative  $\phi$ -entropy.

### 3. ANALYSIS OF THE SEMI-DISCRETE SCHEME

In this section we present the properties of the semi-discrete scheme (16)-(19). Let us state the first main result of this paper which regroups the well-posedness, entropy dissipation and stability properties.

**Theorem 3.1.** *Suppose that the initial data  $f^{\text{in}}$  and the stationary state  $f^\infty$  are positive and consider the semi-discrete scheme (16)-(19) corresponding to (1). Then,*

- *there exists a unique global solution  $f \in \mathcal{C}^1(\mathbb{R}_+; X_{fb})$ ;*
- *there exists two positive constants  $\underline{I}, \bar{I}$  only depending on the initial data  $f^{\text{in}}$  and the stationary state  $f^\infty$ , such that for all  $t \geq 0$  and  $K \in \mathcal{T}$ ,*

$$\underline{I} \leq f_K(t) \leq \bar{I};$$

- the scheme preserves the stationary state  $f^\infty$  and dissipates every relative  $\phi$ -entropy defined in (21), namely for any  $t \geq 0$ ,

$$(23) \quad \frac{d}{dt} H_\phi + D_\phi \leq 0 \quad \text{and} \quad D_\phi \geq 0,$$

where the dissipation  $D_\phi$  is consistent with  $\mathcal{D}_\phi$ .

The key-point in the proof of Theorem 3.1 is to prove the entropy dissipation estimate for any function  $\phi$ . Then we establish the  $L^\infty$  bound on  $f$ . This is done in Section 3.1 and the rest of the proof is detailed in Section 3.2.

Furthermore, in linear cases  $\eta(s) = s$  we can also prove exponential decay rate of the solution to the discrete equilibrium, using the  $\phi$ -entropy inequalities and a discrete Poincaré-Sobolev inequality.

**Theorem 3.2.** (*Exponential return to equilibrium*) Under the assumptions of Theorem 3.1 and for  $\eta(s) = s$ , let  $f$  be the solution of the semi-discrete scheme (16)-(19). Then for  $\xi > 0$  such that  $d_{K,\sigma} \geq \xi d_\sigma$  for all control volume  $K \in \mathcal{T}$  and edge  $\sigma \in \mathcal{E}_K$ , there exists a rate  $\kappa > 0$  depending on the domain,  $\xi$  and  $f^\infty$  (but not on the discretization) and a constant  $C_{0,\infty}$  depending additionally to the initial data such that

$$H_\phi(t) \leq C_{0,\infty} e^{-\kappa t},$$

for all  $\phi$ -entropy satisfying  $2\left(\phi'''\right)^2 \leq \phi''\phi^{IV}$ . In particular, this implies that the semi-discrete solution goes to equilibrium exponentially fast in time

$$\|f(t) - f^\infty\|_{L^1(\Omega)}^2 \leq C_{0,\infty} e^{-\kappa t}.$$

The proof of Theorem 3.2 is given in Section 3.3.

**3.1. Relative entropy dissipation and stability.** In the following we prove the relative entropy inequalities (23) of Theorem 3.1. The result is stated in the following Proposition.

**Proposition 3.3.** Let  $T \in \mathbb{R}_+ \cup \{\infty\}$  and  $f \in C^1(0, T; X_{fb})$ . Then, for any entropy function  $\phi$ , the dissipations associated to  $f$  defined in (22) satisfy the following properties.

- The numerical dissipation  $C_\phi(f)$  is non-negative.
- The physical dissipation rewrites

$$(24) \quad \begin{aligned} D_\phi(f) &= \sum_{\substack{\sigma \in \mathcal{E}_{int} \\ \sigma = K|L}} \tau_\sigma D_{K,\sigma} h D_{K,\sigma} \phi'(h) \eta(f_\sigma^\infty) \\ &+ \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{ext,K}} \tau_\sigma D_{K,\sigma} h D_{K,\sigma} \phi'(h) \eta(f_\sigma^\infty) \geq 0. \end{aligned}$$

Moreover if  $f$  satisfies the scheme (16)-(19), then,

$$(25) \quad \frac{d}{dt} H_\phi + D_\phi = -C_\phi \leq 0.$$

We proceed in three steps to prove Proposition 3.3. First we prove the entropy equality (25). Then, we show that  $D_\phi$  is consistent with its continuous analogue and we finish by the most important property which is the non-negativity of  $C_\phi$  coming from monotony properties of the convective flux. We recall that at the continuous level, this quantity vanishes. This justifies the denomination *numerical (entropy) dissipation*. To prove the non-negativity of the numerical dissipation we will compare it with  $C_\phi^{M_\phi}$  which is the numerical dissipation of a special centered convective flux that depends on the entropy function  $\phi$ . Let us define the latter here as well as some complementary notation.

**Definition 3.4.** A function  $M : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}$  is called a *mean function* if it satisfies for all  $s, t \in \mathbb{R}_+$ ,

- (1)  $M(s, t) = M(t, s)$
- (2)  $M(s, s) = s$

(3) If  $s < t$ , then  $s < M(s, t) < t$

We also define  $M_\sigma : X_{u^b} \rightarrow \mathbb{R}$  by  $M_\sigma(u) = M(u_K, u_L)$  if  $\sigma = K|L$  and  $M_\sigma(u) = M(u_K, u_\sigma)$  otherwise. This is well defined thanks to the symmetry of  $M$ . For any such function  $M$ , we define the centered convective flux associated to  $M$  by

$$F_{K,\sigma}^M(f) = m(\sigma) U_{K,\sigma} M_\sigma(h),$$

and

$$C_\phi^M = \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \phi'(h_K) F_{K,\sigma}^M.$$

Finally for any entropy generating function  $\phi$ , it is elementary to show that

$$M^\phi(s, t) = \frac{\varphi(s) - \varphi(t)}{\phi'(s) - \phi'(t)},$$

where  $\varphi(s) = s\phi'(s) - \phi(s)$ , defines a continuous mean function. We call it the  $\phi$ -mean.

**Remark 3.5.** Let us note that for the 2-entropy generating function  $\phi_2(s) = (s-1)^2$ , the corresponding  $\varphi$ -mean is the arithmetic mean and therefore  $F_{K,\sigma}^{M_{\phi_2}}$  is a classical centered approximation for the convective flux, namely for  $\sigma \in \mathcal{E}_{int}$

$$F_{K,\sigma}^{M_{\phi_2}} = m(\sigma) U_{K,\sigma} \frac{h_K + h_L}{2}.$$

When choosing the generator of the physical entropy  $\phi_1(s) = s \log(s) - s + 1$ , the corresponding mean function is the logarithmic mean reading  $M_{\phi_1}(s, t) = (s-t)/(\log(s) - \log(t))$ .

We are now ready to prove Proposition 3.3.

*Proof of Proposition 3.3.* First note that a simple computation yields

$$\frac{d}{dt} H_\phi = \sum_{K \in \mathcal{T}} m(K) \frac{df_K}{dt} \phi'(h_K) = - \sum_{K \in \mathcal{T}} \phi'(h_K) \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma} = -(C_\phi + D_\phi).$$

using the definition of the dissipations (22) and the scheme (16)-(19). Then, to prove (24), we use (22) and (19) to get

$$D_\phi(f) = - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma D_{K,\sigma} h \phi'(h_K) \eta(f_\sigma^\infty)$$

and the result stems from a discrete integration by parts.

Now, let us prove the non-negativity of  $C_\phi$ . Let  $M$  be any mean function. We start by integrating  $C_\phi - C_\phi^M$  by parts. It yields

$$C_\phi - C_\phi^M = - \sum_{\substack{\sigma \in \mathcal{E}_{int} \\ \sigma = K|L}} (F_{K,\sigma}^{conv} - F_{K,\sigma}^M) D_{K,\sigma}(\phi'(h)) - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{ext,K}} (F_{K,\sigma}^{conv} - F_{K,\sigma}^M) D_{K,\sigma}(\phi'(h)).$$

Now let us just remark that for any  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_{int}$ ,

$$\begin{aligned} - (F_{K,\sigma}^{conv} - F_{K,\sigma}^M) D_{K,\sigma}(\phi'(h)) &= m(\sigma) U_{K,\sigma}^+(g(h_K, h_K) - g(h_K, h_L))(\phi'(h_L) - \phi'(h_K)) \\ &\quad + m(\sigma) U_{K,\sigma}^+(M_\sigma(h) - h_K)(\phi'(h_L) - \phi'(h_K)) \\ &\quad + m(\sigma) U_{K,\sigma}^-(h_L - M_\sigma(h))(\phi'(h_L) - \phi'(h_K)) \\ &\quad + m(\sigma) U_{K,\sigma}^-(g(h_L, h_K) - g(h_L, h_L))(\phi'(h_L) - \phi'(h_K)). \end{aligned}$$

where we used that  $g(s, s) = s$ . If  $\sigma \in \mathcal{E}_{ext}$ , the same equation holds replacing  $h_L$  with  $h_\sigma$ . Therefore, since  $\phi'$  and  $g(s, \cdot)$  are monotonically non-decreasing functions and  $M_\sigma(h)$  is always between  $h_K$  and  $h_L$  (resp.  $h_\sigma$ ), the above quantity is non-negative. Hence,  $C_\phi \geq C_\phi^M$ .

Finally, a simple computation using two integrations by parts yields

$$\begin{aligned}
C_\phi^{M_\phi} &= - \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} F_{K,\sigma}^{M_\phi} D_{K,\sigma}(\phi'(h)) - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} F_{K,\sigma}^{M_\phi} D_{K,\sigma}(\phi'(h)) \\
&= - \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} F_{K,\sigma}^\infty D_{K,\sigma}(\varphi(h)) - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} F_{K,\sigma}^\infty D_{K,\sigma}(\varphi(h)) \\
&= \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^\infty \varphi(h_K) = 0,
\end{aligned}$$

where we used Equation (15) in the last equality. Thus,  $C_\phi \geq C_\phi^{M_\phi} = 0$ .  $\square$

**Remark 3.6.** Let us note that if we had used the fluxes  $F_{K,\sigma}^{M_\phi}$  instead of  $F_{K,\sigma}^{\text{conv}}$  in our scheme, then it would have given,

$$\frac{d}{dt} H_\phi + D_\phi = 0,$$

which is exactly what we get for the continuous problem. However, the scheme would have been  $\phi$ -dependent. With our upwind discretization, we get the whole class of relative entropy inequalities at the cost of an additional numerical dissipation.

The discrete entropy inequalities (23) constitute a large set of Lyapunov functionals that we may use to derive stability properties of the solution. Let us define precisely the stability property of the scheme. Because of the reformulation of (1) into (5) it corresponds to the  $L^\infty$  stability of  $h$ .

**Definition 3.7.** We say that a solution to the semi-discrete scheme (16) is *stable* on  $[0, T)$  if

$$\forall t \in [0, T), \forall K \in \mathcal{T}, \quad h_K(t) \in J,$$

where  $J = [\min(1, \min_K h_K(0)), \max(1, \max_K h_K(0))]$ .

**Remark 3.8.** Since the stationary state  $f^\infty$  is assumed to be positive (component-wise), we introduce the positive constants

$$m_\infty = \min_{K \in \mathcal{T}} \eta(f_K^\infty), \quad M_\infty = \max_{K \in \mathcal{T}} \eta(f_K^\infty).$$

and we set

$$I = [\eta^{-1}(m_\infty \min J), \eta^{-1}(M_\infty \max J)],$$

where  $J$  is defined in Definition 3.7. Assume that  $f \in \mathcal{C}^1(0, T; X_{fb})$  is a solution to the scheme (16)-(19) that is stable in the sense of Definition 3.7. Then,

$$\forall K \in \mathcal{T}, \forall t \in [0, T), \quad f_K(t) \in I.$$

This means that the stability of  $h$ , positivity of  $f^\infty$  and strict monotonicity of  $\eta$  provide the stability of  $f$ .

**Lemma 3.9.** Assume that  $f \in \mathcal{C}^1(0, T; X_{fb})$  is such that for any entropy function  $\phi$ ,

$$\frac{d}{dt} H_\phi + D_\phi \leq 0,$$

with non-negative dissipations  $D_\phi$ . Then it is stable on  $[0, T)$  in the sense of Definition 3.7.

*Proof.* We first restrict to the case where  $J$  is such that  $\inf J < 1 < \sup J$ . Note that the Bernoulli function defined by  $B(x) = x/(\exp(x) - 1)$  is a strictly convex  $\mathcal{C}^2$  function. Then, for any  $\varepsilon > 0$  and  $u_0 \in \mathbb{R}$ , one readily checks that the functions

$$\phi_{\varepsilon, u_0} : u \longmapsto \phi_{\varepsilon, u_0}(u) = \varepsilon \left[ B\left(\frac{u - u_0}{\varepsilon}\right) - B\left(\frac{1 - u_0}{\varepsilon}\right) \right] + B'\left(\frac{1 - u_0}{\varepsilon}\right) (1 - u),$$

are entropy generating functions. Moreover, when  $\varepsilon$  tends to 0, both converges uniformly on any compact subset of  $\mathbb{R}$  to

$$\phi_{\varepsilon, u_0} \longrightarrow \begin{cases} (\cdot - u_0)^+ & \text{if } u_0 > 1 \\ (\cdot - u_0)^- & \text{if } u_0 < 1 \end{cases}.$$

Now, integrating the entropy inequalities, we get in particular that

$$0 \leq H_{\phi_\varepsilon, u_0}(t) \leq H_{\phi_\varepsilon, u_0}(0).$$

If one takes  $u_0 = \inf J$ , then  $u_0 < 1$  and hence passing to the limit  $\varepsilon \rightarrow 0$  in the latter equation leads to

$$0 \leq \sum_{K \in \mathcal{T}} m(K) \int_{f_K^\infty}^{f_K(t)} \left( \frac{\eta(s)}{\eta(f_K^\infty)} - u_0 \right)^- ds \leq \sum_{K \in \mathcal{T}} m(K) \int_{f_K^\infty}^{f_K(0)} \left( \frac{\eta(s)}{\eta(f_K^\infty)} - u_0 \right)^- ds.$$

Using the definition of  $u_0$  one sees that the integrands of the right-hand side vanish. Therefore, for any  $K \in \mathcal{T}$ ,

$$\int_{f_K^\infty}^{f_K(t)} \left( \frac{\eta(s)}{\eta(f_K^\infty)} - u_0 \right)^- ds = 0,$$

which yields that  $h_K(t) \geq u_0$ . The same argument gives the bound from above. Now if  $\inf J = 1$  or  $\sup J = 1$ , then with the same proof we have that uniformly in  $\varepsilon > 0$ ,  $K \in \mathcal{T}$  and  $t \in [0, T)$ ,  $h_K(t) \in J_\varepsilon = [\inf J - \varepsilon, \sup J + \varepsilon]$ . Therefore  $h_K(t) \in J$ .  $\square$

**3.2. Proof of Theorem 3.1.** The Cauchy-Lipschitz theorem yields the existence and uniqueness of a maximal solution  $f \in C^1(0, T; X_{f^b})$  to the Cauchy problem (16)-(19) for some positive final time  $T$ . On this time interval we can apply Proposition 3.3 to get the entropy inequalities (23) for all generating function  $\phi$ . Using Lemma 3.9 and Remark 3.8, we get the existence of finite constants  $\bar{I} = \sup I$  and  $\underline{I} = \inf I$  depending only on the stationary state and the initial data such that

$$\forall t \in [0, T), \forall K \in \mathcal{T}, \quad \underline{I} \leq f_K(t) \leq \bar{I}.$$

Since the solution of the Cauchy problem do not blow up at time  $T$ , it is actually global.

The preservation of the stationary state stems from the fact that  $F_{K,\sigma}^{\text{diss}}(f^\infty) = 0$  for all  $\sigma \in \mathcal{E}_K$  and

$$\sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{\text{conv}}(f^\infty) = \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (U_{K,\sigma}^+ - U_{K,\sigma}^-) = \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^\infty = 0,$$

using (15).

**3.3. Long-time behavior: proof of Theorem 3.2.** In this section, we study the long-time behavior of the discrete solution in the linear case  $\eta(s) = s$ . After providing some preliminary comments and results we prove Theorem 3.2. Our strategy for getting exponential decay to equilibrium from entropy-entropy dissipation properties (23) is the use of a discrete Poincaré-Sobolev type inequality, namely an inequality yielding

$$H_\phi \leq \lambda_{PS} D_\phi,$$

for some constant  $\lambda_{PS}$  only depending on the domain,  $f^\infty$  (hence implicitly on the boundary conditions) and not on the discretization. While the question of the existence of a general  $\phi$ -Poincaré-Sobolev for any entropy generating function  $\phi$ , even in the continuous setting, goes way beyond the scope of this paper (see [6, 3]), it is however possible get such a functional inequality for the particular entropy function  $\phi_2(s) = (s - 1)^2$ . Now, let us detail the whole procedure.

We restrict our class of entropy generating functions to that introduced by Arnold, Markowich, Toscani and Unterreiter in [3], that is those satisfying  $\phi \in C^4(\mathbb{R}_+)$  and

$$\left( \phi''' \right)^2 \leq \frac{1}{2} \phi'' \phi^{IV}.$$

Let us note that the physical and  $p$ -entropies are generated by these entropy functions. The goal of this restriction is to use the consequences of [3, Lemma 2.6] which yield that any such  $\phi$  is bounded from above by a quadratic entropy function, namely

$$\frac{\phi(s)}{\phi''(1)} \leq \phi_2(s) := (s - 1)^2.$$

Therefore, the same inequality holds for the corresponding relative entropies and it suffices to show that the 2-entropy goes to zero exponentially fast in time to get the same result for the whole class relative entropies. The dissipation of 2-entropy is closely related to the discrete  $H^1$  semi-norm

$$|u|_{1,2,\mathcal{T}}^2 = \sum_{\sigma \in \mathcal{E}} \tau_\sigma |D_{K,\sigma} u|^2,$$

for which M. Bessemoulin-Chatard, C. Chainais-Hillairet and F. Filbet proved a discrete Poincaré-Sobolev inequality in [4, Theorem 6], that we shall recall here. Let us first define the  $L^p$  norm

$$\|u\|_{0,p}^p = \sum_{K \in \mathcal{T}} m(K) |u_K|^p.$$

**Proposition 3.10** ([4]). *Suppose that the mesh satisfy the following regularity constraint: there exists  $\xi > 0$  such that  $d_{K,\sigma} \geq \xi d_\sigma$  for all control volume  $K \in \mathcal{T}$  and edge  $\sigma \in \mathcal{E}_K$ . Then there exists a constant  $C$  only depending on the domain such that for all  $u \in X_0$ , it holds*

$$\|u\|_{0,2} \leq \frac{C}{\xi^{1/2}} |u|_{1,2,\mathcal{T}}$$

We are now equipped to prove the long-time behavior result of Theorem 3.2.

*Proof of Theorem 3.2.* Since  $h - 1 \in X_0$ , we may use the expression (24) of the 2-entropy dissipation  $D_{\phi_2}$  in Proposition 3.3 and the result of Proposition 3.10 to get

$$\begin{aligned} D_{\phi_2}(f) &\geq 2 \sum_{\sigma \in \mathcal{E}} \tau_\sigma |D_{K,\sigma}(h-1)|^2 \eta(f_\sigma^\infty) \\ &\geq \frac{2m_\infty \xi^{1/2}}{C} \|h-1\|_{0,2}^2. \end{aligned}$$

Then we notice that  $M_\infty \|h-1\|_{0,2}$  controls  $H_{\phi_2}$  and inject everything in the entropy inequality to get

$$\frac{d}{dt} H_{\phi_2} + \frac{2m_\infty \xi^{1/2}}{C M_\infty} H_{\phi_2} \leq 0,$$

which yields the first result. For the estimate in  $L^1$  it suffices to apply the Hölder inequality to get

$$\|f - f^\infty\|_{0,1} \leq \|f^\infty\|_{0,1}^{1/2} \|f/\sqrt{f^\infty} - \sqrt{f^\infty}\|_{0,2} = \|f^\infty\|_{0,1}^{1/2} H_\varphi^{1/2}.$$

□

#### 4. ANALYSIS OF FULLY DISCRETE SCHEMES

Here, we introduce the explicit and implicit Euler time discretization of the semi-discrete scheme (16)-(19). We denote by  $f^n$  the approximation of  $f$  at time  $t^n = n\Delta t$  and to ease the notation we shall use the superscript  $n$  for any other quantity depending on  $f^n$ .

**4.1. Implicit Euler.** The fully discrete *implicit* scheme is given for all  $K \in \mathcal{T}$  and  $n \in \mathbb{N}$  by

$$(26) \quad \begin{cases} m(K) \frac{f_K^{n+1} - f_K^n}{\Delta t} + \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{n+1} = 0, \\ F_{K,\sigma}^{n+1} = F_{K,\sigma}^{\text{conv}}(f^{n+1}) + F_{K,\sigma}^{\text{diss}}(f^{n+1}), \\ f_K^0 = \frac{1}{m(K)} \int_K f^{\text{in}}(\mathbf{x}) \, d\mathbf{x}, \end{cases}$$

and the fluxes are defined in (18) and (19).

**Theorem 4.1** (Implicit Euler). *Suppose that the initial data  $f^{\text{in}}$  and the stationary state  $f^\infty$  are positive and consider the **fully-discrete implicit scheme** defined by (26),(18) and (19). Then,*

- there exists a unique solution  $f : \mathbb{N} \mapsto X_{fb}$ ;

- there exists two positive constants  $\underline{I}, \bar{I}$  only depending on the initial data  $f^{\text{in}}$  and the stationary state  $f^\infty$ , such that for all  $n \in \mathbb{N}$  and  $K \in \mathcal{T}$ ,

$$\underline{I} \leq f_K^n \leq \bar{I};$$

- the scheme preserves the stationary state  $f^\infty$  and dissipates every relative  $\phi$ -entropy defined in (21), namely for any  $n \in \mathbb{N}$ ,

$$(27) \quad \frac{H_\phi^{n+1} - H_\phi^n}{\Delta t} + D_\phi^{n+1} \leq 0 \quad \text{and} \quad D_\phi^{n+1} \geq 0,$$

where the dissipation  $D_\phi^{n+1}$  is consistent with  $\mathcal{D}_\phi(t^{n+1})$ .

*Proof.* The existence of a unique solution to the implicit scheme can be shown with a fixed point strategy close to that in [16, Remark 4.9] and we do not detail this part. Let us derive the entropy inequality. The Taylor-Young theorem provide the existence of  $\theta_K^{n,n+1} \in (\min(f_K^n, f_K^{n+1}), \max(f_K^n, f_K^{n+1}))$  such that

$$\begin{aligned} e_{\phi,K}^{n+1} - e_{\phi,K}^n &= \int_{f_K^n}^{f_K^{n+1}} \phi' \left( \frac{\eta(s)}{\eta(f_K^\infty)} \right) ds \\ &= \phi'(h_K^{n+1})(f_K^{n+1} - f_K^n) - \frac{1}{2} \psi_K(\theta_K^{n,n+1})(f_K^{n+1} - f_K^n)^2 \\ &= -\frac{\Delta t}{m(K)} \phi'(h_K^{n+1}) \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{n+1} - \frac{1}{2} \psi_K(\theta_K^{n,n+1})(f_K^{n+1} - f_K^n)^2, \end{aligned}$$

where  $\psi_K$  is given by

$$(28) \quad \psi_K : x \mapsto \frac{\eta'(x)}{\eta(f_K^\infty)} \phi'' \left( \frac{\eta(x)}{\eta(f_K^\infty)} \right).$$

Note that  $\psi_K$  is a positive function thanks to the positive monotony of  $\eta$  and  $\phi'$ . With equation (22) this yields

$$\frac{H_\phi^{n+1} - H_\phi^n}{\Delta t} + D_\phi^{n+1} \leq -C_\phi^{n+1} - \frac{1}{2\Delta t} \sum_{K \in \mathcal{T}} \psi_K(\theta_K^{n,n+1})(f_K^{n+1} - f_K^n)^2 m(K) \leq 0,$$

and the rest follows.  $\square$

**4.2. Explicit Euler.** The fully discrete *explicit* scheme is given for all  $K \in \mathcal{T}$  and  $n \in \mathbb{N}$  by

$$(29) \quad \begin{cases} m(K) \frac{f_K^{n+1} - f_K^n}{\Delta t} + \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^n = 0, \\ F_{K,\sigma}^n = F_{K,\sigma}^{\text{conv}}(f^n) + F_{K,\sigma}^{\text{diss}}(f^n), \\ f_K^0 = \frac{1}{m(K)} \int_K f^{\text{in}}(\mathbf{x}) d\mathbf{x}, \end{cases}$$

and the fluxes are defined in (18) and (19). Before stating the result on this scheme, let us introduce

$$(30) \quad a_{K,\sigma} = \begin{cases} d_\sigma \left[ U_{K,\sigma}^+(h_K - g(h_K, h_L)) + U_{K,\sigma}^-(g(h_L, h_K) - h_K) \right] / D_{K,\sigma} h & \text{if } \sigma = K|L, \\ d_\sigma \left[ U_{K,\sigma}^+(h_K - g(h_K, h_\sigma)) + U_{K,\sigma}^-(g(h_\sigma, h_K) - h_K) \right] / D_{K,\sigma} h & \text{otherwise,} \end{cases}$$

with the convention  $a_{K,\sigma} = 0$  if  $D_{K,\sigma} = 0$ . Then mark that we can use this  $a_{K,\sigma}$  to reformulate the convective part of the scheme as a "diffusive term", thank to the incompressibility of  $U_{K,\sigma}$ . Indeed, for all  $K \in \mathcal{T}$ , we have

$$-\sum_{\sigma \in \mathcal{E}_K} \tau_\sigma a_{K,\sigma} D_{K,\sigma} h = \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{\text{conv}} - h_K \sum_{\sigma \in \mathcal{E}_K} m(\sigma) U_{K,\sigma} = \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{\text{conv}},$$

where we used (15). Moreover, because of the monotonicity properties of  $g$ , its local Lipschitz continuity and the fact that  $s = g(s, s)$ , we have that, if  $h \in J$  (cf. Definition 3.7), then

$$0 \leq a_{K,\sigma} \leq C_g U_\infty,$$

where  $C_g$  is the Lipschitz constant of  $g$  on  $J^2$  and  $U_\infty = \max_{K \in \mathcal{T}} \max_{\sigma \in \mathcal{E}_K} d_\sigma |U_{K,\sigma}|$ .

**Theorem 4.2** (Explicit Euler). *Let  $f : \mathbb{N} \mapsto \in X_{fb}$  be defined by the **fully-discrete explicit scheme** (29), (18) and (19). Suppose that the initial data  $f^{in}$  and the stationary state  $f^\infty$  are positive. Then,*

- *there exists a positive constant  $C_{\infty,in}$  depending only on the initial data and the stationary state such that under the CFL condition*

$$\max_{K \in \mathcal{T}} \max_{\sigma \in \mathcal{E}_K} \frac{\tau_\sigma \Delta t}{m(K)} \leq C_{\infty,in},$$

*there exists two positive constants  $\underline{I}, \bar{I}$  only depending on the initial data  $f^{in}$  and the stationary state  $f^\infty$ , such that for all  $n \in \mathbb{N}$  and  $K \in \mathcal{T}$ ,*

$$\underline{I} \leq f_K^n \leq \bar{I};$$

- *if  $\Phi$  is a family of entropy functions with second derivate bounded between  $m_\Phi$  and  $M_\Phi$ , then there exists a positive constant  $\tilde{C}_{\infty,in}$  depending only on the initial data and the stationary state such that for every  $\varepsilon \in (0, 1)$ , under the CFL condition*

$$\max_{K \in \mathcal{T}} \max_{\sigma \in \mathcal{E}_K} \frac{\tau_\sigma \Delta t}{m(K)} \leq \tilde{C}_{\infty,in} \frac{m_\Phi}{M_\Phi} \varepsilon,$$

*$f$  dissipates every relative  $\phi$ -entropy with  $\phi \in \Phi$ , namely for any  $n \in \mathbb{N}$ ,*

$$(31) \quad \frac{H_\phi^{n+1} - H_\phi^n}{\Delta t} + (1 - \varepsilon) D_\phi^n \leq 0 \quad \text{and} \quad D_\phi^n \geq 0$$

*where the dissipation  $D_\phi^n$  is consistent with  $\mathcal{D}_\phi(t^n)$ .*

*Moreover the scheme preserves the stationary state  $f^\infty$ .*

*Proof.* We proceed in three steps. First we derive the entropy-entropy production equality which differ mainly differ from the implicit case by the sign of the remainder term. Then we prove the  $L^\infty$  stability of the scheme in order to achieve the third step, which is the control of the remainder term.

*Entropy equality:* We proceed exactly as in the proof of Theorem 4.1 to get the existence of  $\theta_K^{n,n+1}$  such that

$$(32) \quad \frac{H_\phi^{n+1} - H_\phi^n}{\Delta t} + D_\phi^n \leq -C_\phi^n + \frac{1}{2\Delta t} \sum_{K \in \mathcal{T}} \psi_K(\theta_K^{n,n+1})(f_K^{n+1} - f_K^n)^2 m(K),$$

for  $\psi_K$  defined by (28). Note that the sign of the last term has changed compared to the implicit scheme.

*$L^\infty$  stability:* Let  $\delta > 0$  be such that  $\inf I > \delta > 0$  and define  $M_\eta^\delta = \sup_{s \in I_\delta} \eta'(s)$ , for  $I_\delta := [\inf I - \delta, \sup I + \delta]$ . We want to show by induction on  $n \in \mathbb{N}$  that if

$$(33) \quad \max_{K \in \mathcal{T}} \max_{\sigma \in \mathcal{E}_K} \frac{\tau_\sigma \Delta t}{m(K)} \leq \min \left( \frac{\delta}{N_{\text{edge}} (\sup J - \inf J) (C_g U_\infty + M_\infty)}, \frac{m_\infty}{N_{\text{edge}} M_\eta^\delta (C_g U_\infty + M_\infty)} \right),$$

then  $h_K^n \in J$  for all  $K \in \mathcal{T}$ . In condition (33), the left term of the right-hand side will be useful to prove a first  $L^\infty$  bound that is worse than the induction hypothesis by a margin of  $\delta$ . Then we use the convexity property and the right term of the right-hand side of hypothesis (33) to improve the estimate. By the definition of  $J$  the induction property holds for  $n = 0$ . Suppose that it holds for  $n \in \mathbb{N}$ . The scheme (29) rewrites

$$f_K^{n+1} = f_K^n + \frac{\Delta t}{m(K)} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma (a_{K,\sigma} + \eta(f_K^\infty)) D_{K,\sigma} h^n,$$

and therefore, since (33) (left constant) is satisfied,  $f_K^{n+1} \in I_\delta$  for all  $K \in \mathcal{T}$ . By the mean value theorem, there exists  $g_K^n \in I_\delta$  such that

$$h_K^{n+1} - h_K^n = \frac{\eta'(g_K^n)}{\eta(f_K^\infty)} (f_K^{n+1} - f_K^n).$$



The scheme can then be rewritten as

$$h_K^{n+1} = \left(1 - \frac{\eta'(g_K^n) \Delta t}{\eta(f_K^\infty) m(K)} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma(a_{K,\sigma} + \eta(f_K^\infty))\right) h_K^n + \frac{\eta'(g_K^n) \Delta t}{\eta(f_K^\infty) m(K)} \sum_{\substack{\sigma \in \mathcal{E}_{\text{int},K} \\ \sigma=K|L}} \tau_\sigma(a_{K,\sigma} + \eta(f_K^\infty)) h_L^n \\ + \frac{\eta'(g_K^n) \Delta t}{\eta(f_K^\infty) m(K)} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} \tau_\sigma(a_{K,\sigma} + \eta(f_K^\infty)) h_\sigma^n,$$

which, because of (33) (right constant), provides  $h_K^{n+1}$  as a convex combination of elements of the convex  $J$  and hence  $\{h_K^{n+1}\}_{K \in \mathcal{T}} \subseteq J$ . The CFL constant  $C_{\infty, \text{in}}$  can then be taken as the supremum of the right hand side of (33) when  $\delta \in (0, \inf I)$ .

*Control of the remainder term:*

Using the scheme, the last term in (32) can be estimated with the Cauchy-Schwartz inequality as

$$\frac{\Delta t}{2} \sum_{K \in \mathcal{T}} \frac{1}{m(K)} \psi_K(\theta_K^{n,n+1}) \left( \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^n \right)^2 \leq \frac{N_{\text{edge}} \Delta t M_\Phi M_\eta^0 (C_g U_\infty + M_\infty)^2}{m_\infty} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \frac{\tau_\sigma^2}{m(K)} (D_{K,\sigma} h^n)^2,$$

Then, using that

$$D_\phi^n \geq \frac{m_\Phi m_\infty}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma (D_{K,\sigma} h^n)^2,$$

yields (31) provided that the CFL condition is satisfied with constant

$$\tilde{C}_{\infty, \text{in}} = \frac{m_\infty^2}{2 N_{\text{edge}} M_\eta^0 (C_g U_\infty + M_\infty)^2}.$$

□

**Remark 4.3.** *The CFL constants of the proof seem to depend on the size of the mesh through  $U_\infty$ ,  $m_\infty$ ,  $M_\infty$  and  $N_{\text{edge}}$ . However, if we restrict the class of meshes to those which cannot exceed a certain number of edges by control volume, then the three constants related to the stationary state can be made independent of the mesh assuming consistency of the stationary flux. Of course, even if we do not enter into details on the approximation of the stationary solution, the transient scheme relies on it, therefore consistency is a natural hypothesis.*

## 5. NUMERICAL SIMULATIONS

**5.1. Convergence and order of accuracy.** In this part, we provide a numerical experiment showing the spatial accuracy of our scheme, especially in the long-time dynamic. The test case is the linear drift-diffusion equation with  $\eta(s) = s$  set in one dimension on the domain  $\Omega = (0, 1)$ . We choose the time step is  $\Delta t = 10^{-6}$ , the final time is  $T = 2$  and the time discretization is implicit. Furthermore the advection field  $E(x) = 1$  and the boundary conditions are  $f(t, 0) = 2$  and  $f(t, 1) = 1 + \exp(1)$ . With these parameters the following function

$$f(t, x) = 1 + \exp(x) + \exp\left(\frac{x}{2} - \left(\pi^2 + \frac{1}{4}\right)t\right) \sin(\pi x)$$

is a solution of (1) and converges to the stationary state  $f^\infty(x) = 1 + \exp(x)$  as time goes to infinity.

In order to illustrate the advantage of our approach compared to the one consisting in a direct approximation of (1), we perform numerical simulations using our entropy preserving scheme (16)-(19) and a finite volume scheme applied on  $f$  with an upwinding for the convective terms and centered approximation of the gradient for the diffusive ones.

In Table 1, we measure the  $L^1$  and  $L^\infty$  error between the reconstruction of the approximate solution  $f_N$  obtained on the regular mesh  $(x_i = \Delta x/2 + i\Delta x)_{i \in \{0, \dots, N-1\}}$  of size  $\Delta x = 1/N$  for both schemes.

The error and experimental order are respectively given by

$$e_N^p = \sup_{t \in [0, T]} \|\Pi_N f - f_N\|_{L^p(\Omega)}, \quad k_{2N}^p = |\log(e_{2N}) - \log(e_N)| / \log(2),$$

where  $\Pi_N$  is the projection operator on the mesh, namely

$$\Pi_N f = \sum_{i=0}^{N-1} f_i \mathbb{1}_{[x_{i-1/2}, x_{i+1/2})}$$

with  $f_i$  a numerical approximation of the average of  $f$  on  $(x_{i-1/2}, x_{i+1/2})$ .

$N$	Error $e_N^1$ (16) – (19)	Order	Error $e_N^1$ <b>Upwind</b>	Order	Error $e_N^\infty$ (16) – (19)	Order	Error $e_N^\infty$ <b>Upwind</b>	Order
20	$2.206 \cdot 10^{-3}$		$4.283 \cdot 10^{-3}$		$3.444 \cdot 10^{-3}$		$7.394 \cdot 10^{-3}$	
40	$1.247 \cdot 10^{-3}$	0.82	$2.370 \cdot 10^{-3}$	0.85	$1.959 \cdot 10^{-3}$	0.81	$4.042 \cdot 10^{-3}$	0.87
80	$6.589 \cdot 10^{-4}$	0.92	$1.243 \cdot 10^{-3}$	0.93	$1.038 \cdot 10^{-3}$	0.92	$2.106 \cdot 10^{-3}$	0.94
160	$3.376 \cdot 10^{-4}$	0.96	$6.358 \cdot 10^{-4}$	0.97	$5.329 \cdot 10^{-4}$	0.96	$1.073 \cdot 10^{-3}$	0.97
320	$1.703 \cdot 10^{-4}$	0.99	$3.21 \cdot 10^{-4}$	0.99	$2.690 \cdot 10^{-4}$	0.99	$5.41 \cdot 10^{-4}$	0.99
640	$8.494 \cdot 10^{-5}$	1	$1.606 \cdot 10^{-4}$	1	$1.342 \cdot 10^{-4}$	1	$2.705 \cdot 10^{-4}$	1
1280	$4.184 \cdot 10^{-5}$	1.02	$7.981 \cdot 10^{-5}$	1.01	$6.613 \cdot 10^{-5}$	1.02	$1.344 \cdot 10^{-4}$	1.01

TABLE 1. Experimental spatial order of convergence in  $L^1$  and  $L^\infty$  on  $[0, T)$ .

Both schemes are first order accurate, but we observe that our entropy preserving scheme (16)-(19) is performing better since the numerical error is twice smaller than the classical upwind scheme. Furthermore in Figure 1, we observe that for large time the numerical error corresponding to the entropy preserving scheme (16)-(19) decays to zero and the order of magnitude of the error becomes up to one hundred times smaller than the one corresponding to the upwinding scheme. This numerical test illustrates on a simple example the advantage of the entropy preserving scheme (16)-(19).

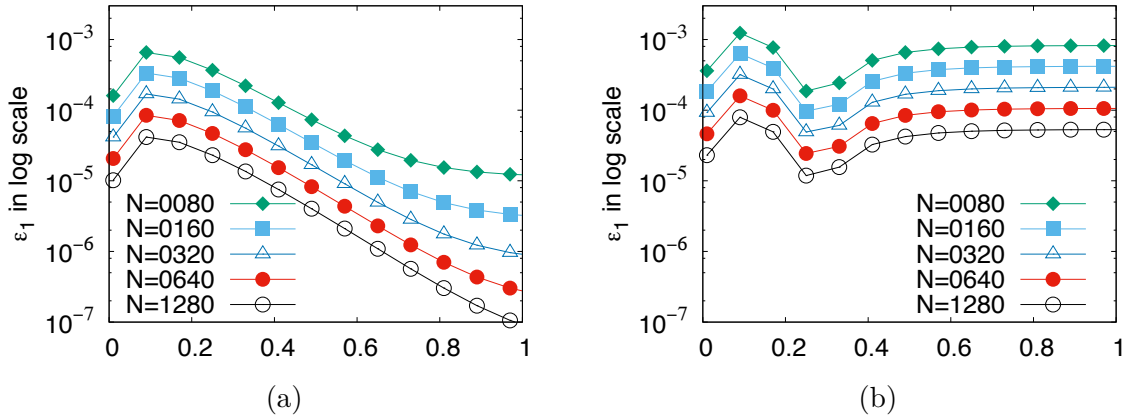


FIGURE 1. **Convergence and order of accuracy.** Time evolution of the  $l_1$  error norm for (a) the entropy preserving scheme and (b) the classical upwind scheme.

**5.2. Fokker-Planck with magnetic field.** We now consider the homogeneous Fokker-Planck equation with an external magnetic field

$$\begin{cases} \frac{\partial f}{\partial t} + \mathbf{v} \wedge \mathbf{B} \cdot \nabla_{\mathbf{v}} f = \nabla_{\mathbf{v}} \cdot (\mathbf{v} f + \nabla_{\mathbf{v}} f) & \text{in } \mathbb{R}^+ \times \mathbb{R}^3, \\ f(t=0) = f_0 & \text{in } \mathbb{R}^3, \end{cases}$$

where the external magnetic field is  $\mathbf{B} = (0, 0, 4)$  whereas the initial datum  $f_0$  is given by the sum of two Gaussian distributions

$$f_0(\mathbf{v}) = \frac{1}{(2\pi)^{3/2}} \left[ \alpha \exp\left(-\frac{|\mathbf{v} - \mathbf{v}_1|^2}{2}\right) + (1 - \alpha) \exp\left(-\frac{|\mathbf{v} - \mathbf{v}_2|^2}{2}\right) \right],$$

with  $\alpha = 3/4$ ,  $\mathbf{v}_1 = (-1, 2, 0)$  and  $\mathbf{v}_2 = (2, -1, 0)$ .

This equation is solved numerically in a bounded domain  $\Omega = (-8, 8)^3$  on various meshes from  $N = 24^3$ , to  $N = 80^3$  points and  $\Delta t = 0.01$  using a time implicit scheme. We choose non homogeneous Dirichlet boundary conditions  $f^b = f^\infty$ , where  $f^\infty$  is the steady state, that is, the Maxwellian distribution

$$f^\infty(\mathbf{v}) = \frac{1}{(2\pi)^{3/2}} \exp\left(-\frac{|\mathbf{v}|^2}{2}\right).$$

Here the knowledge of the steady state  $f^\infty$  allows to compute the stationary flux  $F_{K,\sigma}^\infty$  from a quadrature formula and  $(f_K^\infty)_{K \in \mathcal{T}}$ .

Then for  $h = f/f^\infty$  we define the relative entropy  $H_1(h)$  by

$$H_1(h) := \int_{\mathbb{R}^3} (h - 1)^2 f^\infty d\mathbf{v}$$

and its corresponding dissipation  $D_1(h)$  as

$$D_1(h) := 2 \int_{\mathbb{R}^3} |\nabla h|^2 f^\infty d\mathbf{v}.$$

Finally for  $h = f/f^\infty$  we also define  $H_2(h)$  and  $D_2(h)$  by

$$H_2(h) := \int_{\mathbb{R}^3} [h \log(h) - h + 1] f^\infty d\mathbf{v} \quad \text{and} \quad D_2(h) := \int_{\mathbb{R}^3} \frac{1}{h} |\nabla h|^2 f^\infty d\mathbf{v}.$$

In Figure 2, we represent the time evolution of the entropy  $H_1(h)$ , the dissipation  $D_1(h)$  and the numerical dissipation due to the convective term  $C_1(h)$  in log scale. On the one hand, when  $N \geq 64^3$  points, the entropy and the physical dissipation are well approximated compared to a reference solution : both of them are decreasing function of time and converge to zero with an exponential decay rate. These figures illustrate the convergence to equilibrium exponentially fast. On the other hand, the numerical dissipation of the convective term  $C_1(h)$  converges to zero when the space step goes to zero, but since the scheme is only first order accurate, it is relatively slow. Fortunately, the numerical dissipation  $C_1(h)$  also converges to zero when times goes to infinity exponentially fast, hence it does not affect the accuracy on the decay rate for large time. From these numerical experiments, we get some numerical evidence of the uniform accuracy of the scheme with respect to time.

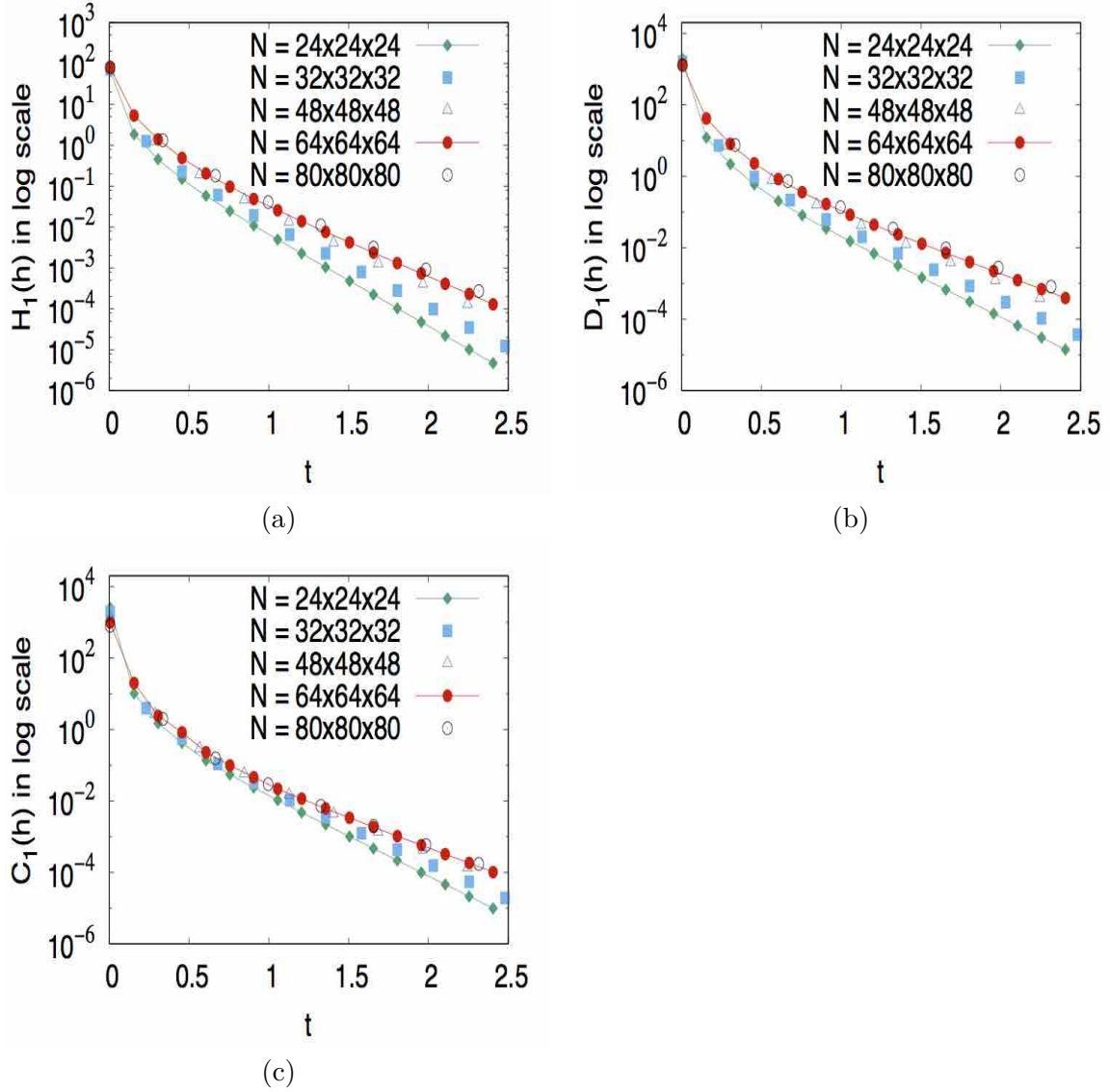


FIGURE 2. **Fokker-Planck equation with an external magnetic field.** Time evolution of (a) the entropy  $H_1(h)$  (b) the physical dissipation  $D_1(h)$  and (c) the numerical dissipation  $C_1(h)$  for different meshes

Furthermore, in Figures 3 and 4 we present a comparison between the physical dissipation  $D_\alpha$ , for  $\alpha \in \{1, 2\}$  which is consistent with (6) and the numerical dissipation  $C_\alpha$  due to the convective (22) for  $N = 24^3$  and  $N = 64^3$  points. First for  $N = 24^3$ , the numerical dissipation  $C_1(h)$  is too large and the decay to equilibrium is amplified. However, for  $N = 64^3$ , the initial dissipation due to the numerical error is smaller than the physical dissipation and then  $C_1(h)$  converges to zero as fast as the physical dissipation  $D_1(h)$ , hence the decay rate obtained for the numerical solution is still consistent with the solution to the Fokker-Planck equation for large time.

**Remark 5.1.** *For the same mesh if we take a larger magnetic field, the numerical dissipation is larger than the physical one. Even if both of them seem to converge to zero with the same decay rate, the dissipation is amplified.*

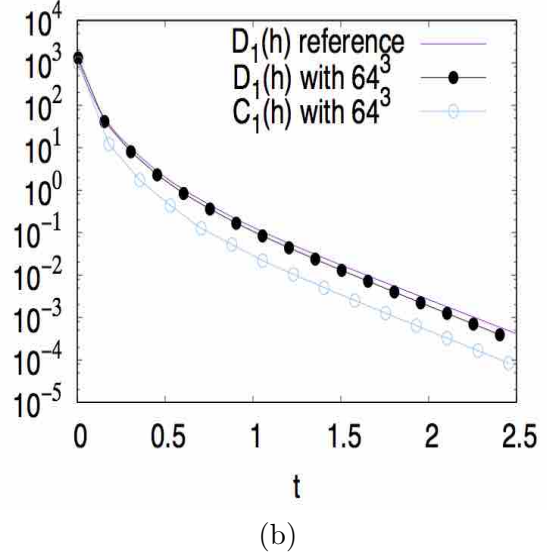
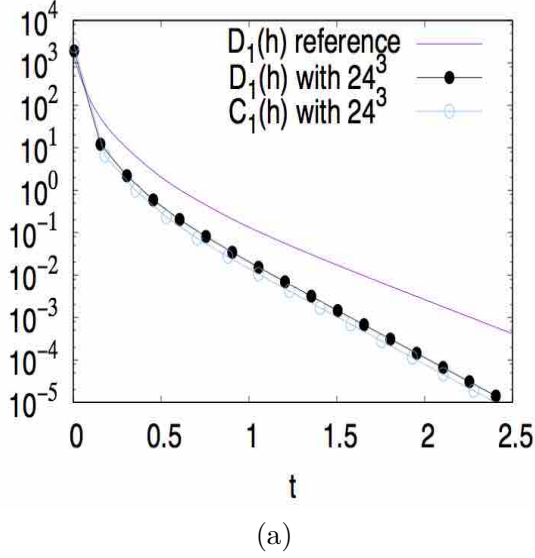


FIGURE 3. **Fokker-Planck equation with an external magnetic field.** Time evolution of the physical dissipation and the numerical dissipation ( $D_1(h), C_1(h)$ ) corresponding to the entropy  $H_1(h)$  with (a)  $N = 24^3$  mesh points (b)  $N = 64^3$  mesh points.

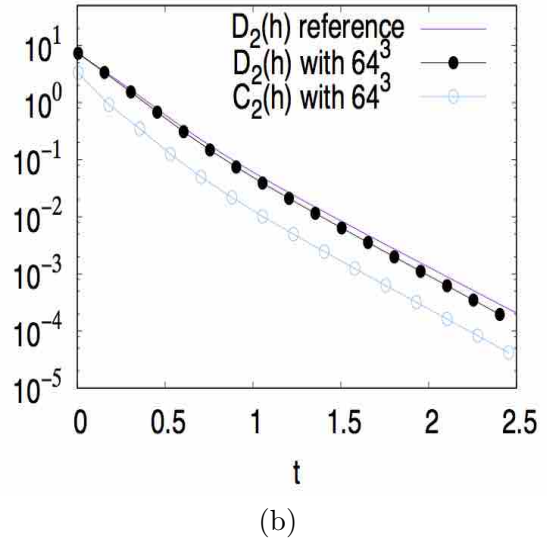
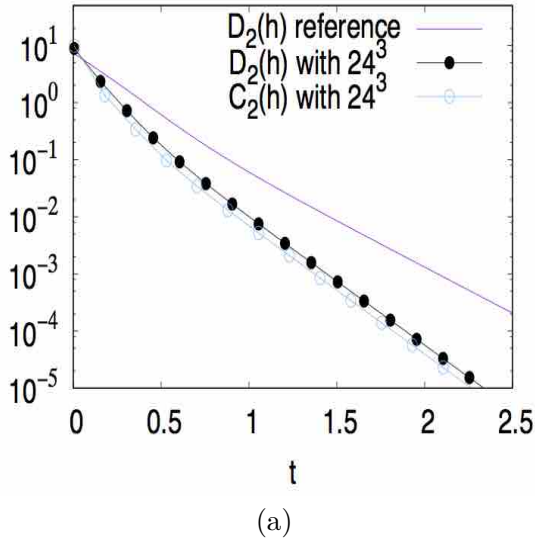


FIGURE 4. **Fokker-Planck equation with an external magnetic field.** Time evolution of the physical dissipation and the numerical dissipation ( $D_2(h), C_2(h)$ ) corresponding to the entropy  $H_2(h)$  with (a)  $N = 24^3$  mesh points (b)  $N = 64^3$  mesh points.

Finally, in Figures 5 and 6, we propose the time evolution of the distribution function at different times. The first column represents an isovalue  $f(t, \mathbf{v}) \equiv 0.01$  of the distribution function whereas the second column is a two dimensional projection in the plane  $v_x - v_y$ , we first observe the effect of the magnetic fields where the two bumps rotate and then under the effect of the Fokker-Planck operator the solution converges to a Maxwellian distribution represented by a sphere in  $\mathbb{R}^3$ .

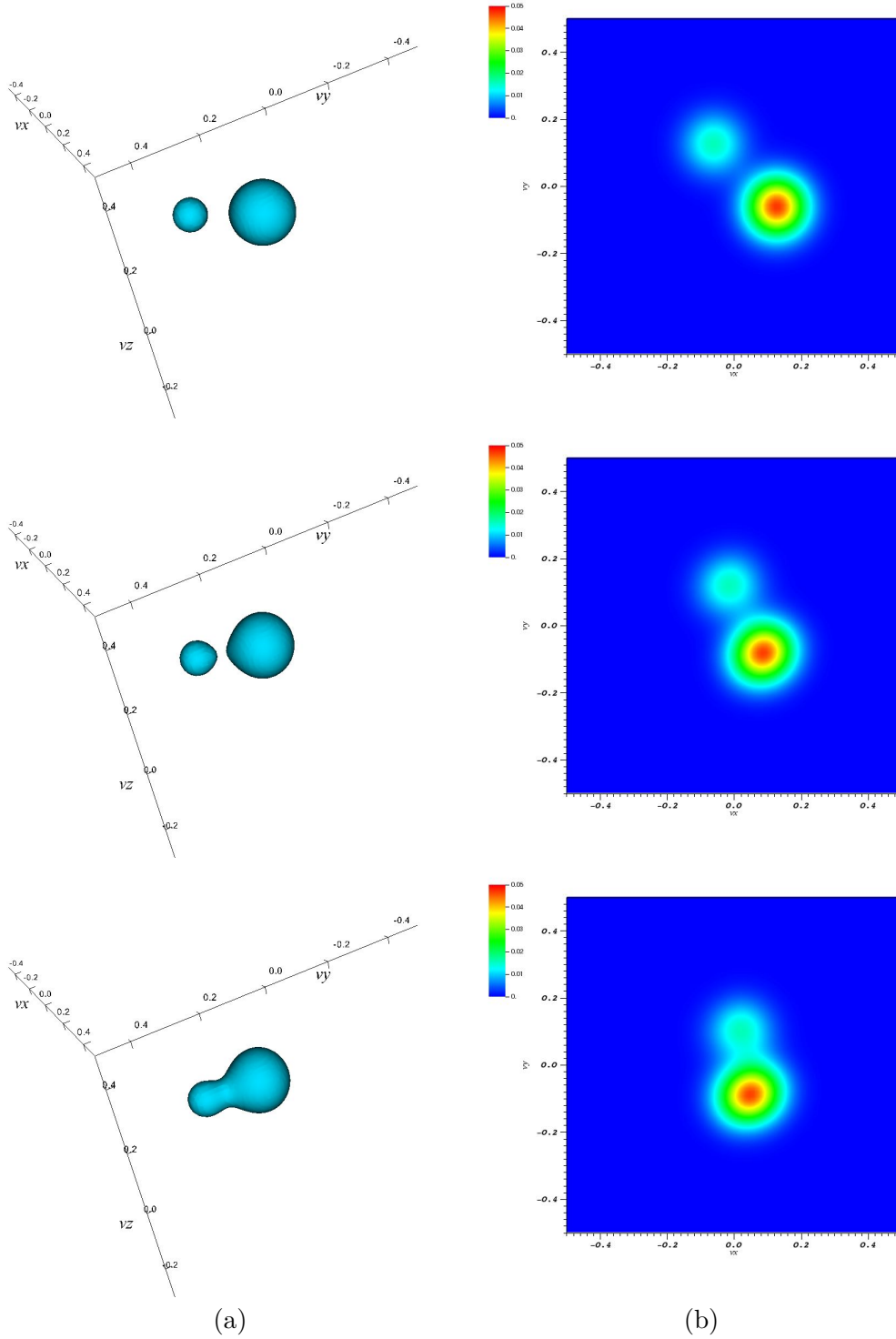


FIGURE 5. Fokker-Planck equation with an external magnetic field. (a) one isovalue (b)  $v_x - v_y$  projection of the distribution in the velocity space at time  $t = 0$ ,  $t = 0.1$  and  $t = 0.2$ .

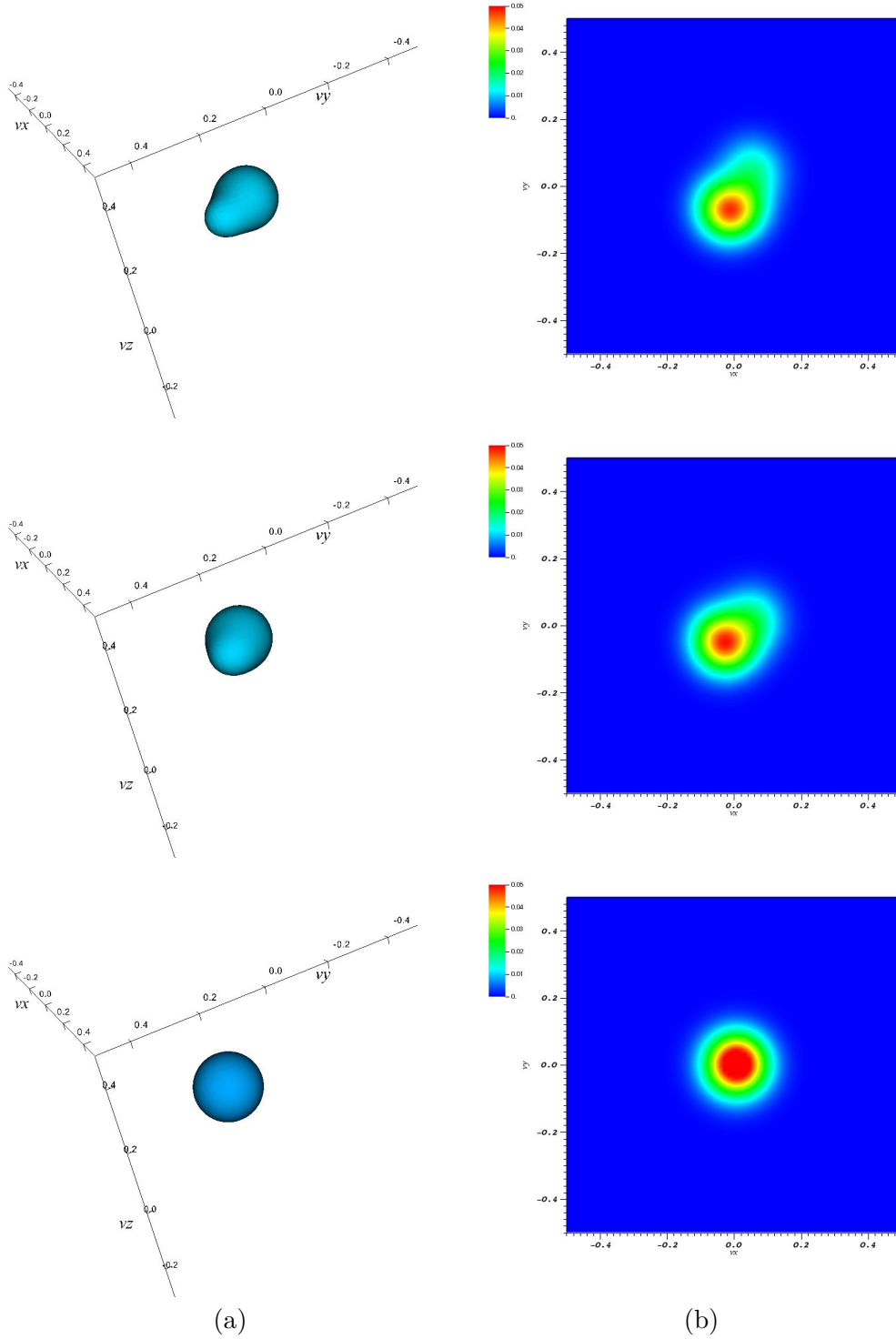


FIGURE 6. **Fokker-Planck equation with an external magnetic field.** (a) one iso-value  $f(t, \mathbf{v}) = 0.01$  (b)  $v_x - v_y$  projection of the distribution in the velocity space at time  $t = 0.3$ ,  $t = 0.4$  and  $t = 0.9$ .

**5.3. Polymer flow in a dilute solution.** We investigate the numerical approximation to the kinetic Fokker-Planck equation for polymers [23]

$$\begin{cases} \frac{\partial F}{\partial t} = -\nabla_{\mathbf{k}} \cdot \left[ \left( \mathbf{A} \mathbf{k} - \frac{1}{2} \nabla_{\mathbf{k}} \Pi(\mathbf{k}) \right) F - \frac{1}{2} \nabla_{\mathbf{k}} F \right], \\ F(t = 0) = F_0 \quad \text{in } \Omega \subset \mathbb{R}^3, \end{cases}$$

where the matrix  $\mathbf{A}$  represents the gradient of an external velocity field and is given by

$$\mathbf{A} = \begin{pmatrix} 1/4 & -1/2 & 0 \\ 1/2 & -1/4 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

and the domain is  $\Omega = (-4, 4)^3$ ,  $\Pi(\mathbf{k}) = |\mathbf{k}|^2/2$ . The initial datum  $F_0$  is given by the sum of two Gaussian distributions

$$F_0(\mathbf{k}) = \frac{1}{2(2\pi)^{3/2}} \left[ \exp\left(-\frac{|\mathbf{k} - \mathbf{k}_1|^2}{2}\right) + \exp\left(-\frac{|\mathbf{k} - \mathbf{k}_2|^2}{2}\right) \right],$$

with  $\mathbf{k}_1 = (-3/2, 1, 0)$  and  $\mathbf{k}_2 = (1, -3/2, 0)$ . This equation is supplemented with homogeneous Neumann boundary conditions such that global mass is conserved. For the numerical simulations we choose various meshes from  $N = 24^3$  to  $64^3$  points with  $\Delta t = 0.01$  using a time implicit scheme. In this case, the steady state is not known, hence the steady state equation is first solved numerically to compute a consistent approximation of the equilibrium  $(f_K^\infty)_{K \in \mathcal{T}}$  and the stationary flux  $F_{K,\sigma}^\infty$ .

Then for  $h = F/f^\infty$  we define the relative entropy  $H_1(h)$  by

$$H_1(h) := \int_{\mathbb{R}^3} (h - 1)^2 f^\infty d\mathbf{k}$$

and its corresponding dissipation  $D_1(h)$  as

$$D_1(h) := 2 \int_{\mathbb{R}^3} |\nabla h|^2 f^\infty d\mathbf{k}.$$

In Figure 7, we represent the time evolution of the entropy  $H_1(h)$ , its dissipation  $D_1(h)$  and the numerical dissipation due to the convective term  $C_1(h)$  in log scale. First when  $N \geq 32^3$  points, the entropy and the physical dissipation are well approximated compared to a reference solution computed with a fine mesh. Once again both of them are decreasing function of time and converge to zero with an exponential decay rate. The numerical dissipation of the convective term  $C_1(h)$  is much smaller than the physical one and also converges to zero when times goes to infinity exponentially fast, hence it does not affect the accuracy on the decay rate for large time.

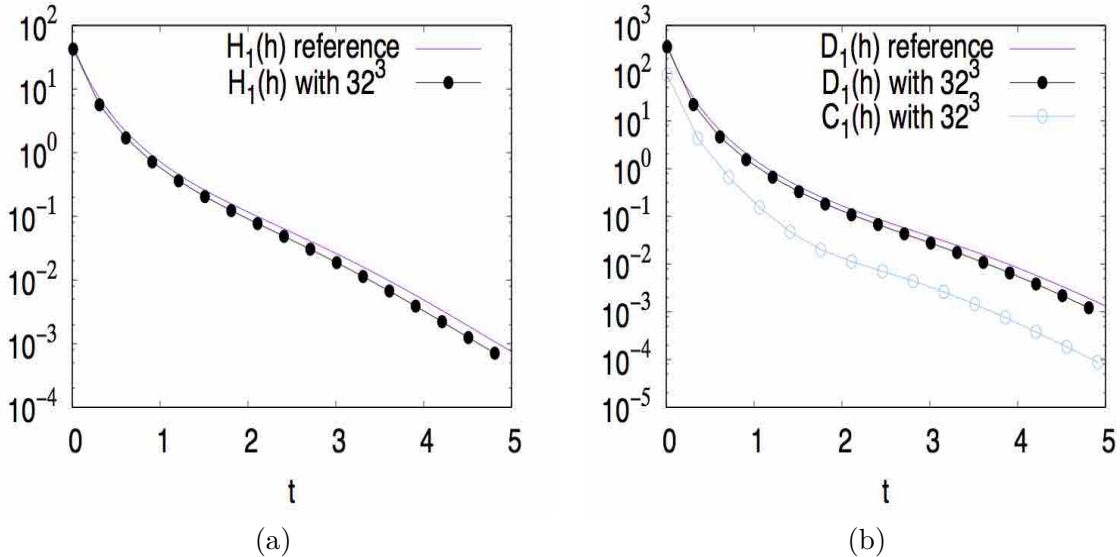


FIGURE 7. **Polymer flow in a dilute solution.** Time evolution of the to the entropy  $H_1(h)$  and the corresponding physical dissipation and numerical dissipation ( $D_1(h), C_1(h)$ ) with  $N = 32^3$  mesh points.

Finally, in Figure 8, we set forth the time evolution of the distribution function at different time. The first column represents an isovalue  $f(t, \mathbf{k}) \equiv 0.02$  of the distribution function whereas the second



column is a two dimensional projection in the plane  $k_x - k_y$ , the solution converges to the discrete steady state, which is consistent to the equilibrium.

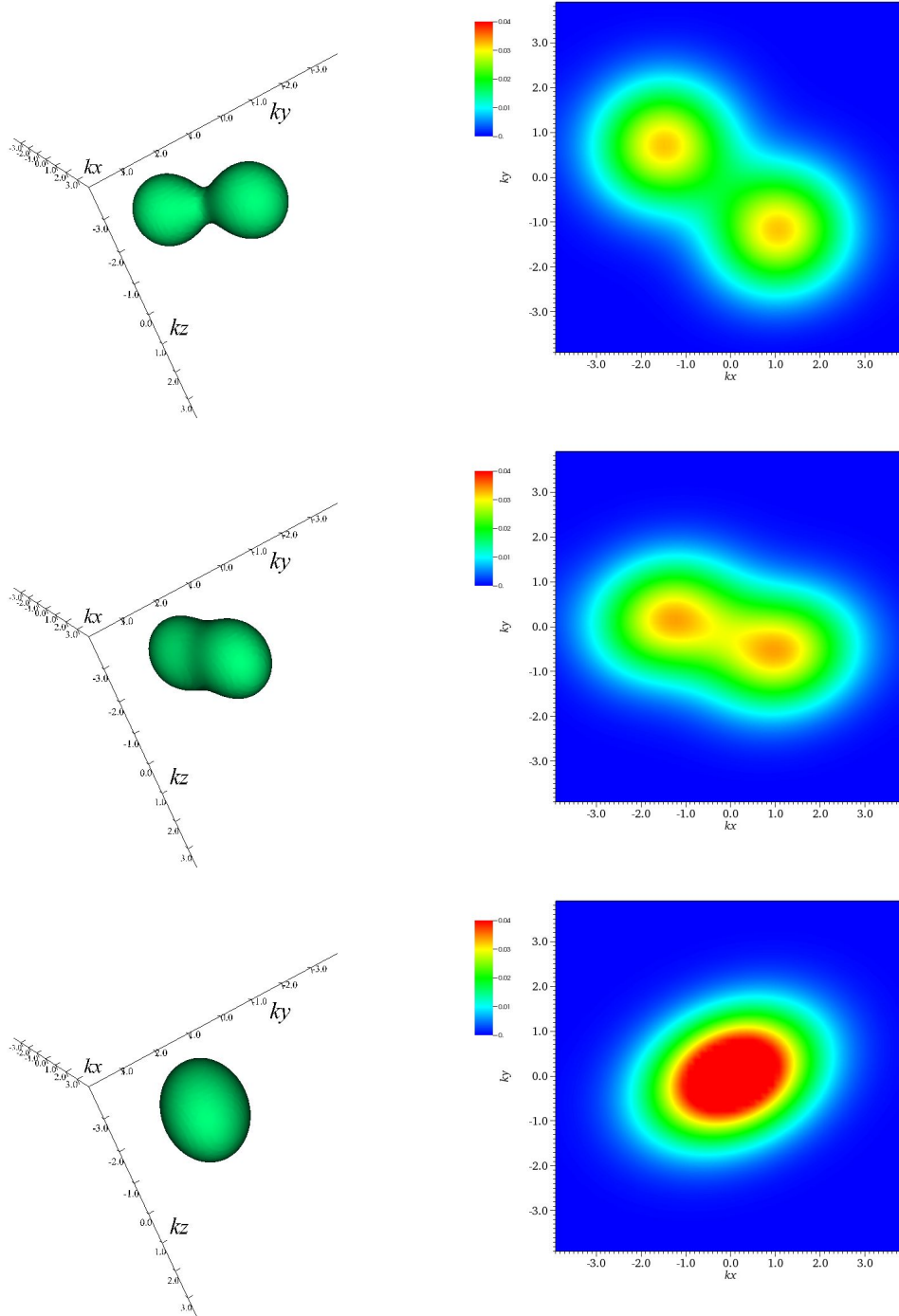


FIGURE 8. **Polymer flow in a dilute solution.** (a) one isovalue  $F(t, \mathbf{k}) = 0.02$  (b)  $k_x - k_y$  projection of the distribution in the  $\mathbf{k}$  space at time  $t = 0.2$ ,  $t = 0.7$  and  $t = 5$ .

5.4. **Porous medium equation.** We finally study the numerical approximation of the porous medium equation

$$\begin{cases} \frac{\partial f}{\partial t} = \Delta f^m \\ f(t=0) = f_0 \quad \text{in } \Omega = (0, 1) \times (-1, 1)^2, \end{cases}$$

with  $m = 2$  and  $f_0 \equiv 1$  together with Dirichlet boundary conditions

$$f^b = \begin{cases} 6 & \text{if } x = 1 \text{ and } y^2 + z^2 \leq 1/8. \\ 1 & \text{else.} \end{cases}$$

This model is nonlinear and without convective terms. Since the steady state is not known, we first compute a numerical approximation  $(f_K^\infty)_{K \in \mathcal{T}}$  and the corresponding stationary flux  $F_{K,\sigma}^\infty$ .

For the numerical simulations we choose various meshes from  $N = 30^3$  to  $60^3$  using a first order time explicit scheme, hence the time step now satisfies a CFL condition  $\Delta t = O(\Delta x^2)$ . For instance as a reference solution, we choose  $N = 60^3$  and  $\Delta t = 0.0001$ . In Figure 9, we represent the time evolution of the relative entropy and the physical and numerical dissipation for  $N = 30^3$ . These results are in good agreement with those obtained using a finer mesh and the numerical dissipation is several order of magnitude smaller than the physical dissipation.

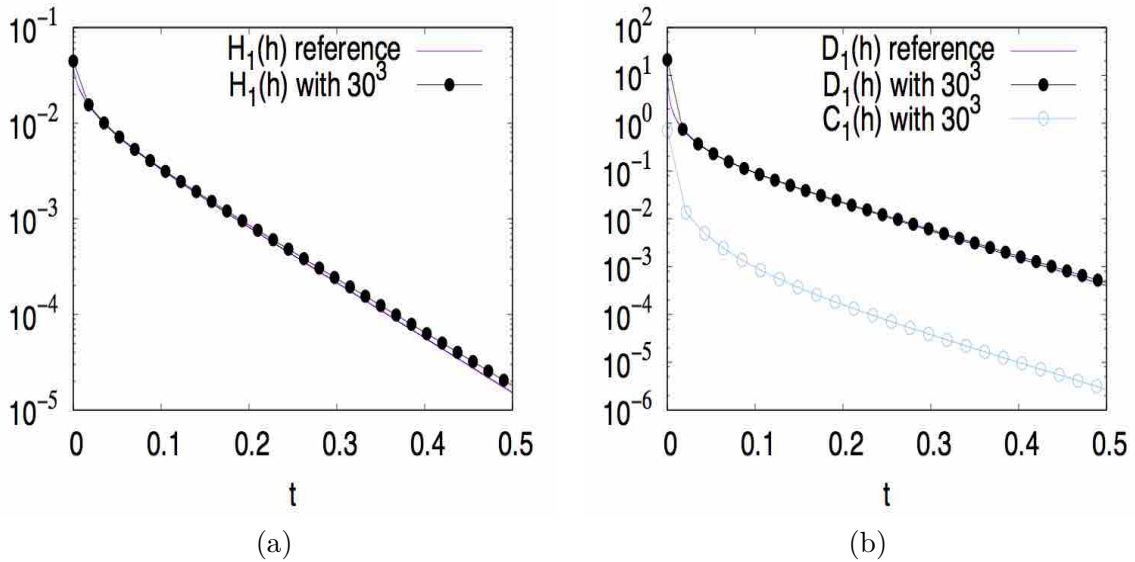


FIGURE 9. **Porous medium equation.** Time evolution of (a) the entropy  $H_1(h)$  (b) the physical dissipation and the numerical dissipation ( $D_1(h)$ ,  $C_1(h)$ ) with  $N = 30^3$  mesh points.

Finally in Figure 10 we represent the time evolution of the distribution function in the three dimensional space. At the initial time the solution is uniform and equation to one, hence the fluid is injected at the boundary  $x = 1$  and we observe how it is diffused in the porous medium at time  $t = 0.05, 0.1$  and  $0.25$ . For  $t \geq 0.5$ , the solution is very close to the equilibrium.

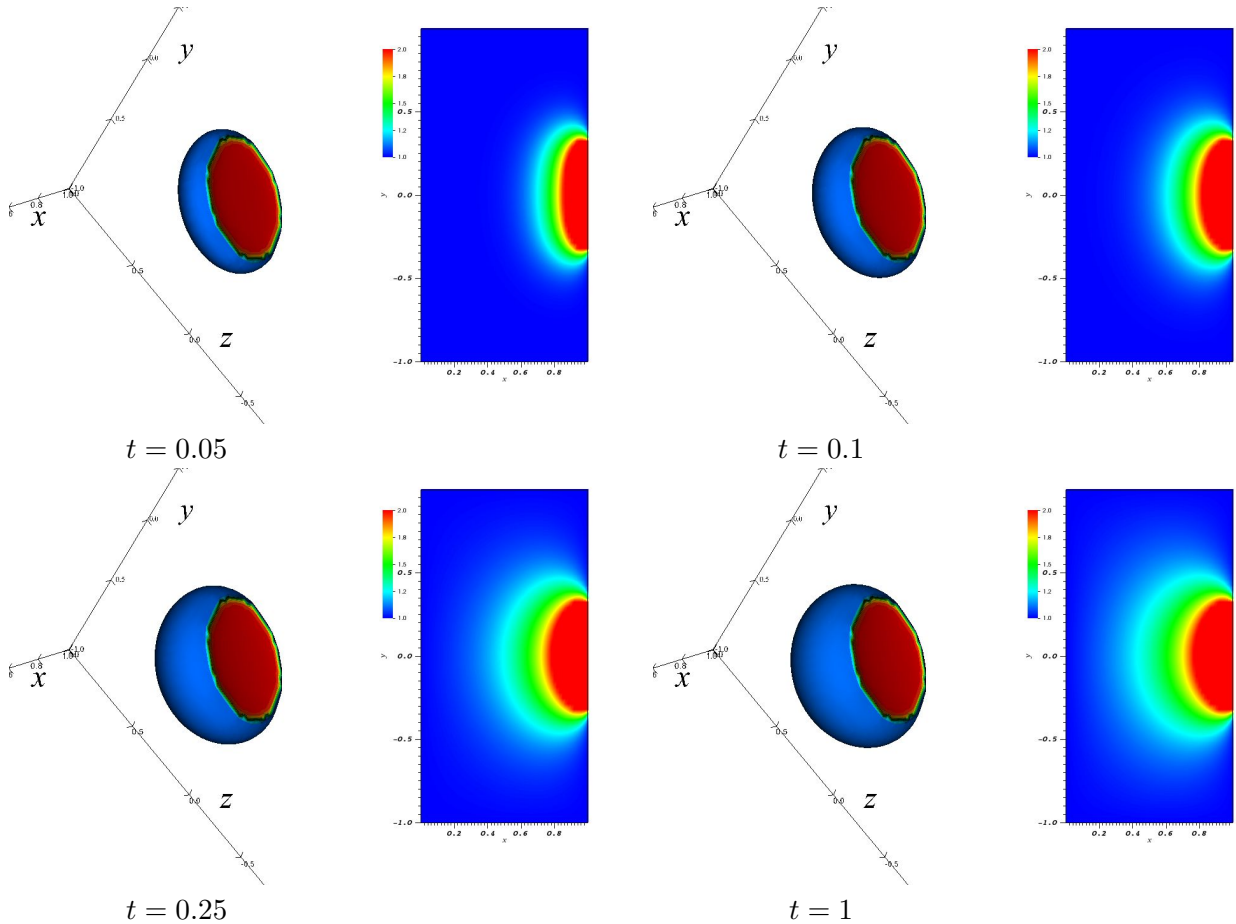


FIGURE 10. **Porous medium equation.** one isovalue  $f(t, \mathbf{x}) = 1.2$  of the distribution and its  $x - y$  projection at time  $t = 0.05$ ,  $t = 0.1$ ,  $t = 0.25$  and  $t = 1$ .

## 6. COMMENTS AND CONCLUSION

In this paper, we have built a scheme for boundary-driven convection-diffusion equations that preserves the relative  $\phi$ -entropy structure of the model. We gave several test cases that confirm the satisfying long time behavior of the numerical scheme in different settings (non-homogeneous Dirichlet/generalized Neumann boundary conditions, explicit and implicit time discretization, linear and nonlinear model).

There are several directions that may be investigated for future work. The first objective is to generalize this scheme to anisotropic diffusions. This mainly requires an adapted discretization of the gradient operator in every direction and there are several papers [17, 14, 15, 18, 10] of Eymard, Herbin, Gallouet, Guichard and Cances are dealing with this issue. Their techniques are based on hybrid finite volume schemes for which the discrete gradient relies on the use of auxiliary unknowns located on the edges between control volumes. This approach seems to be adaptable to our scheme provided that the gradient can be defined in a  $\phi$ -independent way in order to preserve the whole class of Lyapunov functional.

The spirit of our scheme is to start from a consistent discretization of the steady state and build the transient scheme upon the latter to ensure a satisfying behavior in the long-time asymptotic. Therefore, another less specific direction would be the adaptation of this strategy to other types of numerical scheme (Discontinuous Galerkin, Finite elements, *etc.*).

Finally, as we saw in the introduction, many kinetic models (depending on space and velocity variables) such the Vlasov-Fokker-Planck equation or the full dumbbell model for polymers write as the sum of a transport in space and a (convection)-diffusion in velocity. The second part of these models is treated in the present paper. While the diffusion operator is not coercive in all the variables,

thanks to the phase space mixing properties of the transport operator, this still leads to an entropy-diminishing behavior and a trend to a global equilibrium. This property is called hypocoercivity [27] and its preservation by numerical schemes has never been studied to our knowledge and this would be another interesting extension of this work.

ACKNOWLEDGEMENTS. The second author would like to thank Thierry Dumont for his kind help on sparse matrix routines.

## REFERENCES

- [1] Franz Achleitner, Anton Arnold, and Dominik Stürzer. Large-time behavior in non-symmetric fokker-planck equations. 2015.
- [2] Anton Arnold, Eric Carlen, and Qiangchang Ju. Large-time behavior of non-symmetric Fokker-Planck type equations. *Commun. Stoch. Anal.*, 2(1):153–175, 2008.
- [3] Anton Arnold, Peter Markowich, Giuseppe Toscani, and Andreas Unterreiter. On convex Sobolev inequalities and the rate of convergence to equilibrium for Fokker-Planck type equations. *Comm. Partial Differential Equations*, 26(1-2):43–100, 2001.
- [4] Marianne Bessemoulin-Chatard, Claire Chainais-Hillairet, and Francis Filbet. On discrete functional inequalities for some finite volume schemes. *IMA J. Numer. Anal.*, 35(3):1125–1149, 2015.
- [5] Marianne Bessemoulin-Chatard and Francis Filbet. A finite volume scheme for nonlinear degenerate parabolic equations. *SIAM J. Sci. Comput.*, 34(5):B559–B583, 2012.
- [6] Thierry Bodineau, Joel Lebowitz, Clément Mouhot, and Cédric Villani. Lyapunov functionals for boundary-driven nonlinear drift-diffusion equations. *Nonlinearity*, 27(9):2111–2132, 2014.
- [7] F. Bouchut and J. Dolbeault. On long time asymptotics of the Vlasov-Fokker-Planck equation and of the Vlasov-Poisson-Fokker-Planck system with Coulombic and Newtonian potentials. *Differential Integral Equations*, 8(3):487–514, 1995.
- [8] Martin Burger, José A. Carrillo, and Marie-Therese Wolfram. A mixed finite element method for nonlinear diffusion equations. *Kinet. Relat. Models*, 3(1):59–83, 2010.
- [9] Clément Cancès and Cindy Guichard. Numerical analysis of a robust entropy-diminishing Finite Volume scheme for parabolic equations with gradient structure. working paper or preprint, 2015.
- [10] Clément Cancès and Cindy Guichard. Convergence of a nonlinear entropy diminishing control volume finite element scheme for solving anisotropic degenerate parabolic equations. *Math. Comp.*, 85(298):549–580, 2016.
- [11] Claire Chainais-Hillairet and Francis Filbet. Asymptotic behaviour of a finite-volume scheme for the transient drift-diffusion model. *IMA J. Numer. Anal.*, 27(4):689–716, 2007.
- [12] Claire Chainais-Hillairet, Ansgar Jüngel, and Stefan Schuchnigg. Entropy-dissipative discretization of nonlinear diffusion equations and discrete Beckner inequalities. *ESAIM Math. Model. Numer. Anal.*, 50(1):135–162, 2016.
- [13] I. Csiszár. Information-type measures of difference of probability distributions and indirect observations. *Studia Sci. Math. Hungar.*, 2:299–318, 1967.
- [14] R. Eymard, T. Gallouët, and R. Herbin. A cell-centered finite-volume approximation for anisotropic diffusion operators on unstructured meshes in any space dimension. *IMA J. Numer. Anal.*, 26(2):326–353, 2006.
- [15] R. Eymard, T. Gallouët, and R. Herbin. Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes SUSHI: a scheme using stabilization and hybrid interfaces. *IMA J. Numer. Anal.*, 30(4):1009–1043, 2010.
- [16] Robert Eymard, Thierry Gallouët, and Raphaële Herbin. Finite volume methods. In *Handbook of numerical analysis, Vol. VII*, Handb. Numer. Anal., VII, pages 713–1020. North-Holland, Amsterdam, 2000.
- [17] Robert Eymard, Thierry Gallouët, and Raphaële Herbin. A finite volume scheme for anisotropic diffusion problems. *C. R. Math. Acad. Sci. Paris*, 339(4):299–302, 2004.
- [18] Robert Eymard, Cindy Guichard, and Raphaële Herbin. Small-stencil 3D schemes for diffusive flows in porous media. *ESAIM Math. Model. Numer. Anal.*, 46(2):265–290, 2012.
- [19] Francis Filbet and Chi-Wang Shu. Approximation of hyperbolic models for chemosensitive movement. *SIAM J. Sci. Comput.*, 27(3):850–872 (electronic), 2005.
- [20] Maxime Herda. On massless electron limit for a multispecies kinetic system with external magnetic field. *J. Differential Equations*, 260(11):7861–7891, 2016.
- [21] Benjamin Jourdain, Claude Le Bris, Tony Lelièvre, and Félix Otto. Long-time asymptotics of a multiscale model for polymeric fluid flows. *Arch. Ration. Mech. Anal.*, 181(1):97–148, 2006.
- [22] Solomon Kullback. *Information theory and statistics*. John Wiley and Sons, Inc., New York; Chapman and Hall, Ltd., London, 1959.
- [23] Nader Masmoudi. Well-posedness for the FENE dumbbell model of polymeric flows. *Comm. Pure Appl. Math.*, 61(12):1685–1714, 2008.
- [24] Mark A Peletier. Variational modelling: Energies, gradient flows, and large deviations. *arXiv preprint arXiv:1402.1990*, 2014.
- [25] Andreas Unterreiter, Anton Arnold, Peter Markowich, and Giuseppe Toscani. On generalized Csiszár-Kullback inequalities. *Monatsh. Math.*, 131(3):235–253, 2000.

- [26] Juan Luis Vázquez. *The porous medium equation*. Oxford Mathematical Monographs. The Clarendon Press, Oxford University Press, Oxford, 2007. Mathematical theory.
- [27] Cédric Villani. Hypocoercivity. *Mem. Amer. Math. Soc.*, 202(950):iv+141, 2009.

FRANCIS FILBET

INSTITUT DE MATHÉMATIQUES DE TOULOUSE,  
UNIVERSITÉ TOULOUSE III & INSTITUT UNIVERSITAIRE DE FRANCE,  
BÂTIMENT 1R3, 118, ROUTE DE NARBONNE  
F-31062, TOULOUSE CEDEX 9 FRANCE

E-MAIL: FRANCIS.FILBET@MATH.UNIV-TOULOUSE.FR

MAXIME HERDA

INSTITUT CAMILLE JORDAN,  
UNIVERSITÉ CLAUDE BERNARD LYON 1, CNRS UMR 5208,  
43 BLVD. DU 11 NOVEMBRE 1918  
F-69622, VILLEURBANNE CEDEX FRANCE

E-MAIL: HERDA@MATH.UNIV-LYON1.FR