

Structure preserving schemes for Fokker–Planck equations with nonconstant diffusion matrices

*Original*

Structure preserving schemes for Fokker–Planck equations with nonconstant diffusion matrices / Loy, N.; Zanella, M.. - In: MATHEMATICS AND COMPUTERS IN SIMULATION. - ISSN 0378-4754. - 188:(2021), pp. 342-362. [10.1016/j.matcom.2021.04.018]

*Availability:*

This version is available at: 11583/2918598 since: 2021-09-01T15:29:33Z

*Publisher:*

Elsevier

*Published*

DOI:10.1016/j.matcom.2021.04.018

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

Elsevier postprint/Author's Accepted Manuscript

© 2021. This manuscript version is made available under the CC-BY-NC-ND 4.0 license  
<http://creativecommons.org/licenses/by-nc-nd/4.0/>. The final authenticated version is available online at:  
<http://dx.doi.org/10.1016/j.matcom.2021.04.018>

(Article begins on next page)

# Structure preserving schemes for Fokker-Planck equations with nonconstant diffusion matrices

Nadia Loy \*

Mattia Zanella<sup>†</sup>

## Abstract

In this work we consider an extension of a recently proposed structure preserving numerical scheme for nonlinear Fokker-Planck-type equations to the case of nonconstant full diffusion matrices. While in existing works the schemes are formulated in a one-dimensional setting, here we consider exclusively the two-dimensional case. We prove that the proposed schemes preserve fundamental structural properties like nonnegativity of the solution without restriction on the size of the mesh and entropy dissipation. Moreover, all the methods presented here are at least second order accurate in the transient regimes and arbitrarily high order for large times in the hypothesis in which the flux vanishes at the stationary state. Suitable numerical tests will confirm the theoretical results.

**Keywords:** Fokker-Planck equations, positivity preserving, structure preserving methods, finite difference schemes

**Mathematics Subject Classification:** 35Q70, 35Q84, 65N06

## 1 Introduction

We are interested in nonlinear Fokker-Planck equations describing the evolution of a multivariate distribution function  $f(w, t) \geq 0$ , with  $t \geq 0, w \in \Omega \subseteq \mathbb{R}^2$  of the following form

$$\begin{cases} \partial_t f(w, t) = \nabla_w \cdot \mathcal{F}(w, t), & \mathcal{F}(w, t) = \mathcal{B}[f](w, t)f(w, t) + \nabla_w \cdot (\mathbb{D}(w)f(w, t)), & w \in \Omega, \\ \mathcal{F}(w, t) \cdot \mathbf{n}(w) = 0, & w \in \partial\Omega, \\ f(w, 0) = f_0(w), & w \in \Omega, \end{cases} \quad (1)$$

where  $\Omega \subset \mathbb{R}^2$  is bounded and  $\mathbf{n}(w)$  is the outward normal unit vector defined for  $w \in \partial\Omega$ . In particular the no-flux boundary condition  $\mathcal{F}(w, t) \cdot \mathbf{n}(w) = 0, w \in \partial\Omega$  guarantees conservation of mass in  $\Omega$ , i.e.  $\int_{\Omega} f(w, t) dw = \int_{\Omega} f_0(w) dw \forall t \geq 0$ . The drift term  $\mathcal{B}[f](\cdot, t)$  can classically be defined as the following nonlocal bounded operator

$$\begin{aligned} \mathcal{B}[f](\cdot, t) : \Omega &\mapsto \mathbb{R}^2 \\ w &\mapsto \mathcal{B}[f](w, t) = \int_{\Omega} P(w, w_*) (w - w_*) f(w_*, t) dw_*, \end{aligned} \quad (2)$$

where  $P(\cdot, \cdot) : \Omega \times \Omega \rightarrow \mathbb{R}^+$ . Therefore, the drift term  $\mathcal{B}[f](w, t)$  depends on time only through  $f(w, t)$ . In (1) we consider a nonconstant diffusion matrix  $\mathbb{D}(w)$  which is supposed to be symmetric and positive definite in the interior of  $\Omega$ , while it vanishes at the border, i.e.

$$\mathbb{D}(w) = 0, \quad w \in \partial\Omega. \quad (3)$$

---

\*Department of Mathematical Sciences “G. L. Lagrange”, Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129 Torino, Italy, (nadia.loy@polito.it)

<sup>†</sup>Department of Mathematics “F. Casorati”, Via A. Ferrata 5, 27100 Pavia, Italy (mattia.zanella@unipv.it)

Kinetic-type equations with general diffusion often arise in the derivation of aggregate descriptions of many particles systems. Without intending to review the very huge literature on this topic we mention [6, 2, 3, 12, 19] for applications to collective phenomena, [16, 26, 30, 44] for related models in self-organized biological aggregations, and [21, 41, 32, 27, 42, 43] for their relation with Boltzmann-type modelling. Kinetic equations have a strong physical interpretation as they describe the time evolution of probability density functions that describe the statistical distribution of the microscopic variables. Therefore, their solution should be nonnegative. Moreover, their trend to equilibrium is studied through an entropy functional that is dissipated in time and minimized by a unique stationary equilibrium. The necessity to deal with a general diffusion matrix arises from various applications where heterogeneity appears in the evolution of the distribution function. Of course, this gives rise to a genuinely multi-dimensional problem whose stationary state makes the divergence of the flux vanish.

In this manuscript we concentrate on the the two-dimensional problem and on the special case where also the flux of the problem vanishes at the stationary equilibrium. We will develop in this setting finite difference numerical schemes for the introduced problem that preserve structural properties like nonnegativity of the solution, entropy dissipation and that approximate with arbitrary accuracy the stationary state of the problem. Furthermore, the methods here developed are second order accurate in the transient regime and do not require restrictions on the mesh size. The schemes with the introduced features are usually referred to as structure preserving schemes (SP). The methods here derived are based on recent works in this direction [11, 18, 34, 35] and follow pioneering works on linear Fokker-Planck equations [14, 25], see also [8, 9, 20, 38, 39]. We refer to [4, 5, 10, 13, 36] for related methods in the case of degenerate diffusion and to [22] for a recent survey on methods preserving steady states of balance laws.

In more details the paper is organized as follows. In Section 2 we derive the structure preserving scheme. We will compare the obtained scheme with recent results for 1D problems. Hence, in Section 3 we prove nonnegativity of the numerical solution in the case of explicit and semi-implicit time integration. Sufficient conditions will be explicated in terms of bounds on the time step. The trend to equilibrium is then investigated in Section 4 in the case of linear problems, here we prove that the constructed SP scheme dissipates the numerical entropy. Finally in Section 5 we present several applications of the schemes in Fokker-Planck problems describing emerging patterns in interacting systems. Some conclusions are reported at the end of the manuscript.

## 2 Structure preserving schemes and full nonconstant diffusion matrices

In this section we focus on the design of a numerical scheme for the nonlinear Fokker-Planck equation with general diffusion matrix of the form (1).

### 2.1 Stationary states

We first observe that in (1) the two dimensional flux  $\mathcal{F}(w, t) = [\mathcal{F}^x(w, t), \mathcal{F}^y(w, t)]^T$  is given by

$$\mathcal{F}(w, t) = \mathcal{C}(w, t)f(w, t) + \mathbb{D}(w)\nabla_w f(w, t), \quad (4)$$

with  $\mathbb{D}(w)$  a nonconstant diffusion matrix of the form

$$\mathbb{D}(w) = \begin{bmatrix} \mathbb{D}^{1,1}(w) & \mathbb{D}^{1,2}(w) \\ \mathbb{D}^{2,1}(w) & \mathbb{D}^{2,2}(w) \end{bmatrix},$$

such that  $\mathbb{D} \in \mathcal{C}^2(\Omega)$  and is symmetric and positive definite in the interior of  $\Omega$ . Since we considered  $\mathbb{D}(w)$  symmetric and positive definite its determinant is strictly positive, i.e.  $|\mathbb{D}(w)| > 0, \forall w \in \Omega$ .

We remark that, as  $|\mathbb{D}(w)| = \mathbb{D}_{1,1}(w)\mathbb{D}_{2,2}(w) - \mathbb{D}_{1,2}(w)\mathbb{D}_{2,1}(w) = \mathbb{D}_{1,1}(w)\mathbb{D}_{2,2}(w) - \mathbb{D}_{1,2}^2(w)$ , the diagonal elements of the diffusion matrix, i.e.  $\mathbb{D}^{1,1}(w)$  and  $\mathbb{D}^{2,2}(w)$ , cannot vanish  $\forall w \in \Omega$ . Moreover, in (4) we considered  $\mathcal{C}(w, t) = [\mathcal{C}^x(w, t), \mathcal{C}^y(w, t)]^T$  where

$$\mathcal{C}^x(w, t) = \mathcal{B}[f]^x(w, t) + \partial_x \mathbb{D}^{1,1}(w) + \partial_y \mathbb{D}^{2,1}(w),$$

$$\mathcal{C}^y(w, t) = \mathcal{B}[f]^y(w, t) + \partial_x \mathbb{D}^{1,2}(w) + \partial_y \mathbb{D}^{2,2}(w),$$

and  $\mathcal{B}[f](w, t) = [\mathcal{B}[f]^x(w, t), \mathcal{B}[f]^y(w, t)]$ . In (4) we indicated with  $\partial_x, \partial_y$  the partial derivatives in the directions defined by the components  $w_x$  and  $w_y$ , respectively. Therefore the components of the flux  $\mathcal{F}(w, t)$  are given by

$$\mathcal{F}^x(w, t) = \mathcal{C}^x(w, t)f(w, t) + \mathbb{D}^{1,1}(w)\partial_x f(w, t) + \mathbb{D}^{1,2}(w)\partial_y f(w, t), \quad (5)$$

$$\mathcal{F}^y(w, t) = \mathcal{C}^y(w, t)f(w, t) + \mathbb{D}^{2,1}(w)\partial_x f(w, t) + \mathbb{D}^{2,2}(w)\partial_y f(w, t). \quad (6)$$

We want to approximate the stationary state, i.e. the multivariate distribution function  $f^\infty(w)$  satisfying

$$\nabla \cdot \mathcal{F}^\infty(w) = 0, \quad (7)$$

where  $\mathcal{F}^\infty(w) = \mathcal{C}(w)f^\infty(w) + \mathbb{D}(w)\nabla_w f^\infty(w)$ . We remark that now  $\mathcal{C} = \mathcal{C}(w)$  depends only on  $w$  as, in general,  $\mathcal{C}$  depends on time only through  $f$  and we are now considering the stationary state. A sufficient condition for  $f^\infty(w)$  to be a stationary state is that it makes the flux vanish, i.e.

$$\mathcal{F}^\infty(w) = 0 \quad (8)$$

that corresponds to have a constant  $\mathcal{F}^\infty(w)$  that may be zero by fixing accordingly the no-flux boundary condition. Equation (8) is a necessary condition only in one-dimension, while in a two-dimensional setting there may be stationary states satisfying (7) such that  $\mathcal{F}^\infty(w) \neq 0$ . We concentrate on the problems in which (8) is satisfied.

Let us observe that in 2D the condition  $\mathcal{F}^\infty(w) = 0$  defines a decoupled system of Fokker-Planck equations if  $\mathbb{D}^{1,2}(w) = \mathbb{D}^{2,1}(w) = 0$ ,  $w \in \Omega$ , that can be solved through standard schemes. In 1D, several numerical strategies to catch the emerging equilibrium have been designed. Schemes for Fokker-Planck-type equations have been studied in the community: without intending to review the huge literature in this direction we mention schemes for linear drift-diffusion-type problems [9, 14, 25, 38] together with related entropy methods [8, 13], and recent developments for the general energy-decaying problems [2, 34]. In particular, in [34] an analogous nonlinear Fokker-Planck equation with diagonal diffusion matrix was considered. In the present work, we will derive the scheme under the hypothesis in which the flux vanishes at the stationary state.

Let us consider  $\mathbb{D}^{1,1}(w), \mathbb{D}^{2,2}(w) \neq 0 \forall w \in \Omega$ , we can define the following quasi-stationary system for the components of the flux

$$\begin{aligned} \mathbb{D}^{1,1}(w)\partial_x f(w, t) &= -f(w, t)\mathcal{C}^x - \mathbb{D}^{1,2}(w)\partial_y f(w, t), \\ \mathbb{D}^{2,2}(w)\partial_y f(w, t) &= -f(w, t)\mathcal{C}^y - \mathbb{D}^{2,1}(w)\partial_x f(w, t). \end{aligned} \quad (9)$$

The latter system is quasi-stationary because  $\mathcal{C}^x(w, t)$  and  $\mathcal{C}^y(w, t)$  depend on  $f(w, t)$ . Let us observe that, thanks to the introduction of the matrix characterizing the diffusion the equations (9) are not decoupled unless  $\mathbb{D}$  is diagonal. From (9) we have

$$\begin{aligned} \left( \mathbb{D}^{1,1}(w) - \frac{\mathbb{D}^{1,2}(w)\mathbb{D}^{2,1}(w)}{\mathbb{D}^{2,2}(w)} \right) \partial_x f(w, t) &= -f(w, t) \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right), \\ \left( \mathbb{D}^{2,2}(w) - \frac{\mathbb{D}^{1,2}(w)\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \right) \partial_y f(w, t) &= -f(w, t) \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right). \end{aligned} \quad (10)$$

In the following we will adopt the notations

$$\mathcal{D}^1(w) = \mathbb{D}^{1,1}(w) - \frac{\mathbb{D}^{1,2}(w)\mathbb{D}^{2,1}(w)}{\mathbb{D}^{2,2}(w)}, \quad \mathcal{D}^2(w) = \mathbb{D}^{2,2}(w) - \frac{\mathbb{D}^{1,2}(w)\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)}, \quad (11)$$

that are positive quantities since  $\mathbb{D}$  is positive definite in the interior of  $\Omega$ .

It is worth stressing how in the case  $\mathbb{D}^{1,2}(w) = \mathbb{D}^{2,1}(w) = 0$ , the two equations in (10) can be decoupled and we basically recover the classical quasi-stationary formulation in each direction, we refer to [25, 34] for more details on the concept of quasi-equilibrium distribution and to [31] for further applications. Furthermore, we remark how system (10) is in general not *analytically* solvable, except in some special cases due to the nonlinearity on the right hand side and the intrinsically coupled nature of the system. We overcome this difficulty in the quasi steady-state approximation integrating the equations of system (10) over numerical grids.

## 2.2 Derivation of the scheme

Let us consider  $\Omega = [a, b] \times [a, b]$  and let us introduce the regular grid

$$W = \left\{ w_{i,j} = (w_{x,i}, w_{y,j}) \in \Omega \mid w_{x,i} = a + i\Delta w, w_{y,j} = a + j\Delta w, i, j = 0, \dots, N, \Delta w = \frac{b-a}{N} \right\}.$$

We shall also define the mid points grid as

$$W^{\text{mid}} = \left\{ w_{i+1/2, j+1/2} = (w_{x, i+1/2}, w_{y, j+1/2}) \in \Omega \mid \right. \\ \left. w_{x, i+1/2} = a + \frac{i\Delta w}{2}, w_{y, j+1/2} = a + \frac{j\Delta w}{2}, i, j = 0, \dots, N-1 \right\}$$

and we remark that  $W^{\text{mid}}$  is in the interior of  $\Omega$ . Without loss of generality, and to avoid unnecessary complications, we consider a square domain with an uniform grid, i.e. with square cells; anyway what presented in the following can be easily generalized to the case of rectangular cells, in which  $w_{x, i+1} - w_{x, i} = \Delta w_1$  and  $w_{y, j+1} - w_{y, j} = \Delta w_2$ . Next, if we integrate the two equations in (10) with respect to  $w_x$  on  $[w_{i,j}, w_{i+1,j}]$  and with respect to  $w_y$  on  $[w_{i,j}, w_{i,j+1}]$ , respectively, we have

$$\int_{w_{i,j}}^{w_{i+1,j}} \frac{\partial_x f(w, t)}{f(w, t)} dw_x = - \int_{w_{i,j}}^{w_{i+1,j}} \frac{1}{\mathcal{D}^1(w)} \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right) dw_x, \\ \int_{w_{i,j}}^{w_{i,j+1}} \frac{\partial_y f(w, t)}{f(w, t)} dw_y = - \int_{w_{i,j}}^{w_{i,j+1}} \frac{1}{\mathcal{D}^2(w)} \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right) dw_y,$$

leading respectively to

$$\frac{f(w_{i+1,j}, t)}{f(w_{i,j}, t)} = \exp \left\{ - \int_{w_{i,j}}^{w_{i+1,j}} \frac{1}{\mathcal{D}^1(w)} \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right) dw_x \right\} \quad (12)$$

and

$$\frac{f(w_{i,j+1}, t)}{f(w_{i,j}, t)} = \exp \left\{ - \int_{w_{i,j}}^{w_{i,j+1}} \frac{1}{\mathcal{D}^2(w)} \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right) dw_y \right\}. \quad (13)$$

Let us denote  $f_{i,j}(t)$  an approximation of  $f(w_{i,j}, t)$  over the grid  $W$ . Let us introduce the following finite difference scheme where we keep the time continuous

$$\frac{d}{dt} f_{i,j}(t) = \frac{\mathcal{F}_{i+1/2,j}^x(t) - \mathcal{F}_{i-1/2,j}^x(t)}{\Delta w} + \frac{\mathcal{F}_{i,j+1/2}^y(t) - \mathcal{F}_{i,j-1/2}^y(t)}{\Delta w}, \quad (14)$$

where the right hand side is a numerical approximation of the operator  $\nabla_w \cdot \mathcal{F}(w, t)$  on the grid  $W$  at time  $t > 0$ . The quantities  $\mathcal{F}_{i\pm 1/2,j}^x(t)$ ,  $\mathcal{F}_{i,j\pm 1/2}^y(t)$  are the numerical flux functions relative to the introduced numerical discretization. We want to define the numerical fluxes analogously to [34], where they give a second order definition for the two components of the numerical flux, i.e.  $\mathcal{F}_{i+1/2,j}^x(t)$  and  $\mathcal{F}_{i,j+1/2}^y(t)$  are combinations of the grid points  $i+1$  and  $i$ ,  $j+1$  and  $j$  respectively. In the rest of this section we will omit the explicit dependency on time.

In particular, in [34], where  $\mathbb{D}^{1,2}(w) = \mathbb{D}^{2,1}(w) = 0$ , the authors define the numerical fluxes as

$$\begin{aligned}\mathcal{F}_{i+1/2,j}^x(t) &= \tilde{\mathcal{C}}_{i+1/2,j}^x(t) \tilde{f}_{i+1/2,j}(t) + \mathbb{D}_{i+1/2,j}^{1,1} \frac{f_{i+1,j}(t) - f_{i,j}(t)}{\Delta w}, \\ \mathcal{F}_{i,j+1/2}^y(t) &= \tilde{\mathcal{C}}_{i,j+1/2}^y(t) \tilde{f}_{i,j+1/2}(t) + \mathbb{D}_{i,j+1/2}^{2,2} \frac{f_{i,j+1}(t) - f_{i,j}(t)}{\Delta w},\end{aligned}\tag{15}$$

where  $\tilde{f}_{i+1/2,j}(t)$ ,  $\tilde{f}_{i,j+1/2}(t)$  are classically defined as

$$\begin{aligned}\tilde{f}_{i+1/2,j}(t) &= (1 - \delta_{i+1/2,j}(t)) f_{i+1,j}(t) + \delta_{i+1/2,j}(t) f_{i,j}(t), \\ \tilde{f}_{i,j+1/2}(t) &= (1 - \delta_{i,j+1/2}(t)) f_{i,j+1}(t) + \delta_{i,j+1/2}(t) f_{i,j}(t),\end{aligned}\tag{16}$$

see [14, 29, 34]. The weight functions  $\delta_{i+1/2,j}(t)$ ,  $\delta_{i,j+1/2}(t)$  are hence defined in such a way that they have values in  $(0, 1)$  and, thus,  $\tilde{f}_{i+1/2,j}(t)$  and  $\tilde{f}_{i,j+1/2}(t)$  are convex combinations of  $f_{i+1,j}(t)$ ,  $f_{i,j}(t)$  and  $f_{i,j+1}(t)$ ,  $f_{i+1,j}(t)$  respectively.

In the present setting, since the extra diagonal terms of the diffusion matrix do not vanish, the definition of the numerical fluxes must be modified accordingly. In particular, we shall write as an extension of (15)

$$\begin{aligned}\mathcal{F}_{i+1/2,j}^x(t) &= \tilde{\mathcal{C}}_{i+1/2,j}^x(t) \tilde{f}_{i+1/2,j}(t) + \mathbb{D}_{i+1/2,j}^{1,1} \frac{f_{i+1,j}(t) - f_{i,j}(t)}{\Delta w} + \mathbb{D}_{i+1/2,j}^{1,2} [\partial_y f]_{i,j}(t), \\ \mathcal{F}_{i,j+1/2}^y(t) &= \tilde{\mathcal{C}}_{i,j+1/2}^y(t) \tilde{f}_{i,j+1/2}(t) + \mathbb{D}_{i,j+1/2}^{2,2} \frac{f_{i,j+1}(t) - f_{i,j}(t)}{\Delta w} + \mathbb{D}_{i,j+1/2}^{2,1} [\partial_x f]_{i,j}(t),\end{aligned}\tag{17}$$

where  $[\partial_y f]_{i,j}(t)$  and  $[\partial_x f]_{i,j}(t)$  are numerical approximations of the partial derivatives  $\partial_y f(w, t)$  and  $\partial_x f(w, t)$  **that we need to determine**. As we want to perform a directional splitting, we have to determine the approximations  $[\partial_y f]_{i,j}(t)$  and  $[\partial_x f]_{i,j}(t)$  in the complementary direction with respect to the one of the differentiation, i.e. as a combination of  $f_{i+1,j}(t)$ ,  $f_{i,j}(t)$  and  $f_{i,j+1}(t)$ ,  $f_{i,j}(t)$  respectively. In order to obtain such approximations, in addition to  $\mathcal{F}_{i+1/2,j}^x(t) = 0$  and  $\mathcal{F}_{i,j+1/2}^y(t) = 0$ , we consider the discretization of the two components of the numerical fluxes in the complementary direction, i.e. we discretize  $\mathcal{F}^x(w, t)$  in the  $y$  direction and  $\mathcal{F}^y(w, t)$  in the  $x$  direction:

$$\begin{aligned}\mathcal{F}_{i,j+1/2}^x(t) &= \tilde{\mathcal{C}}_{i,j+1/2}^x(t) \tilde{f}_{i,j+1/2}(t) + \mathbb{D}_{i,j+1/2}^{1,2} \frac{f_{i,j+1}(t) - f_{i,j}(t)}{\Delta w} + \mathbb{D}_{i,j+1/2}^{1,1} [\partial_x f]_{i,j}(t), \\ \mathcal{F}_{i+1/2,j}^y(t) &= \tilde{\mathcal{C}}_{i+1/2,j}^y(t) \tilde{f}_{i+1/2,j}(t) + \mathbb{D}_{i+1/2,j}^{2,1} \frac{f_{i+1,j}(t) - f_{i,j}(t)}{\Delta w} + \mathbb{D}_{i+1/2,j}^{2,2} [\partial_y f]_{i,j}(t).\end{aligned}\tag{18}$$

By making the latter vanish, i.e.  $\mathcal{F}_{i,j+1/2}^x(t) = 0$  and  $\mathcal{F}_{i+1/2,j}^y(t) = 0$ , we find the following numerical approximations in  $w_{i,j}$  of the partial derivatives  $[\partial_y f]_{i,j}(t)$  and  $[\partial_x f]_{i,j}(t)$  in the complementary direction with respect to the one of the differentiation

$$[\partial_y f]_{i,j}(t) = -\frac{1}{\mathbb{D}_{i+1/2,j}^{2,2}} \left[ \tilde{\mathcal{C}}_{i+1/2,j}^y(t) \tilde{f}_{i+1/2,j}(t) + \mathbb{D}_{i+1/2,j}^{2,1} \frac{f_{i+1,j}(t) - f_{i,j}(t)}{\Delta w} \right],\tag{19}$$

and

$$[\partial_x f]_{i,j}(t) = -\frac{1}{\mathbb{D}_{i,j+1/2}^{1,1}} \left[ \tilde{\mathcal{C}}_{i,j+1/2}^x(t) \tilde{f}_{i,j+1/2}(t) + \mathbb{D}_{i,j+1/2}^{1,2} \frac{f_{i,j+1}(t) - f_{i,j}(t)}{\Delta w} \right],\tag{20}$$

where  $\tilde{f}_{i+1/2,j}(t)$ ,  $\tilde{f}_{i,j+1/2}(t)$  are given by (16). By substituting (19) and (20) in Eq (17) we obtain

$$\mathcal{F}_{i+1/2,j}^x(t) = \tilde{\mathcal{G}}_{i+1/2,j}^x(t) \tilde{f}_{i+1/2,j}(t) + \mathcal{D}_{i+1/2,j}^1 \frac{f_{i+1,j}(t) - f_{i,j}(t)}{\Delta w},\tag{21a}$$

$$\mathcal{F}_{i,j+1/2}^y(t) = \tilde{\mathcal{G}}_{i,j+1/2}^y(t) \tilde{f}_{i,j+1/2}(t) + \mathcal{D}_{i,j+1/2}^2 \frac{f_{i,j+1}(t) - f_{i,j}(t)}{\Delta w},\tag{21b}$$

where  $\tilde{f}_{i+1/2,j}(t)$ ,  $\tilde{f}_{i,j+1/2}(t)$  are expressed as in (16) and

$$\begin{aligned}\tilde{\mathcal{G}}_{i+1/2,j}^x(t) &= \tilde{\mathcal{C}}_{i+1/2,j}^x(t) - \frac{\mathbb{D}_{i+1/2,j}^{1,2}}{\mathbb{D}_{i+1/2,j}^{2,2}} \tilde{\mathcal{C}}_{i+1/2,j}^y(t), \\ \tilde{\mathcal{G}}_{i,j+1/2}^y(t) &= \tilde{\mathcal{C}}_{i,j+1/2}^y(t) - \frac{\mathbb{D}_{i,j+1/2}^{2,1}}{\mathbb{D}_{i,j+1/2}^{1,1}} \tilde{\mathcal{C}}_{i,j+1/2}^x(t).\end{aligned}\tag{22}$$

We shall now equate to zero the two components of the numerical flux (21). By setting (21a) to zero, where  $\tilde{f}_{i+1/2,j}(t)$  is defined as in (16) and  $\tilde{\mathcal{G}}_{i+1/2,j}^x(t)$  as in (22), we obtain

$$f_{i+1,j}(t)(1 - \delta_{i+1/2,j}(t))\tilde{\mathcal{G}}_{i+1/2,j}^x(t) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} + f_{i,j}(t)\delta_{i+1/2,j}(t)\tilde{\mathcal{G}}_{i+1/2,j}^x(t) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} = 0$$

and, therefore

$$\frac{f_{i+1,j}(t)}{f_{i,j}(t)} = \frac{-\delta_{i+1/2,j}(t)\tilde{\mathcal{G}}_{i+1/2,j}^x(t) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w}}{(1 - \delta_{i+1/2,j}(t))\tilde{\mathcal{G}}_{i+1/2,j}^x(t) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w}}.\tag{23}$$

Analogously, equating (21b) to zero gives

$$\frac{f_{i,j+1}(t)}{f_{i,j}(t)} = \frac{-\delta_{i,j+1/2}(t)\tilde{\mathcal{G}}_{i,j+1/2}^y(t) + \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w}}{(1 - \delta_{i,j+1/2}(t))\tilde{\mathcal{G}}_{i,j+1/2}^y(t) + \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w}},\tag{24}$$

where, as a consequence of the definition (11), we have

$$\begin{aligned}\mathcal{D}_{i+1/2,j}^1 &= \mathbb{D}^{1,1}(w_{i+1/2,j}) - \frac{\mathbb{D}^{1,2}(w_{i+1/2,j})\mathbb{D}^{2,1}(w_{i+1/2,j})}{\mathbb{D}^{2,2}(w_{i+1/2,j})} > 0, \\ \mathcal{D}_{i,j+1/2}^2 &= \mathbb{D}^{2,2}(w_{i,j+1/2}) - \frac{\mathbb{D}^{1,2}(w_{i,j+1/2})\mathbb{D}^{2,1}(w_{i,j+1/2})}{\mathbb{D}^{1,1}(w_{i,j+1/2})} > 0.\end{aligned}$$

We now need to define suitable weight functions  $\delta_{i+1/2,j}(t)$ ,  $\delta_{i,j+1/2}(t)$  and numerical drifts  $\tilde{\mathcal{C}}^x(w, t)$ ,  $\tilde{\mathcal{C}}^y(w, t)$  so that the method preserves the steady state of the problem with arbitrary accuracy and so that its numerical solution defines nonnegative solutions without additional restrictions on the grid  $\Delta w$ . By equating analytical and the numerical form of the flux, i.e.  $f(w_{i+1,j}, t)/f(w_{i,j}, t)$  in (12) with  $f_{i+1,j}(t)/f_{i,j}(t)$  in (23), and  $f(w_{i,j+1}, t)/f(w_{i,j}, t)$  in (13) with  $f_{i,j+1}(t)/f_{i,j}(t)$  in (24), and setting

$$\begin{aligned}\tilde{\mathcal{C}}_{i+1/2,j}^x(t) &= \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \int_{w_{i,j}}^{w_{i+1,j}} \frac{\mathcal{C}^x(w, t)}{\mathcal{D}^1(w)} dw_x, \\ \tilde{\mathcal{C}}_{i+1/2,j}^y(t) &= \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \int_{w_{i,j}}^{w_{i+1,j}} \frac{\mathcal{C}^y(w, t)}{\mathcal{D}^1(w)} dw_x,\end{aligned}$$

and

$$\begin{aligned}\tilde{\mathcal{C}}_{i,j+1/2}^x(t) &= \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \int_{w_{i,j}}^{w_{i,j+1}} \frac{\mathcal{C}^x(w, t)}{\mathcal{D}^2(w)} dw_y, \\ \tilde{\mathcal{C}}_{i,j+1/2}^y(t) &= \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \int_{w_{i,j}}^{w_{i,j+1}} \frac{\mathcal{C}^y(w, t)}{\mathcal{D}^2(w)} dw_y,\end{aligned}$$

we finally get

$$\delta_{i+1/2,j}(t) = \frac{1}{\lambda_{i+1/2,j}(t)} + \frac{1}{1 - \exp(\lambda_{i+1/2,j}(t))}, \quad \delta_{i,j+1/2}(t) = \frac{1}{\lambda_{i,j+1/2}(t)} + \frac{1}{1 - \exp(\lambda_{i,j+1/2}(t))}, \quad (25)$$

where

$$\begin{aligned} \lambda_{i+1/2,j}(t) &= \int_{w_{i,j}}^{w_{i+1,j}} \frac{1}{\mathcal{D}^1(w)} \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right) dw_x = \frac{\Delta w}{\mathcal{D}_{i+1/2,j}^1} \tilde{\mathcal{G}}_{i+1/2,j}^x(t), \\ \lambda_{i,j+1/2}(t) &= \int_{w_{i,j}}^{w_{i,j+1}} \frac{1}{\mathcal{D}^2(w)} \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right) dw_y = \frac{\Delta w}{\mathcal{D}_{i,j+1/2}^2} \tilde{\mathcal{G}}_{i,j+1/2}^y(t). \end{aligned} \quad (26)$$

We have the following result

**Theorem 1.** *The numerical flux defined by (17) with (19)-(20) is given by (21) with  $\tilde{\mathcal{G}}_{i+1/2,j}^x(t)$ ,  $\tilde{\mathcal{G}}_{i,j+1/2}^y(t)$  defined in (22) and with  $\delta_{i+1/2,j}(t)$ ,  $\delta_{i,j+1/2}(t)$  defined in (25). The numerical flux (21) vanishes when the flux (5)-(6) *vanishes* over the cell  $[w_{i,j}, w_{i+1,j}] \times [w_{i,j}, w_{i,j+1}]$ . The nonlinear weights defined in (25)-(26) are such that  $\delta_{i\pm 1/2,j}(t) \in (0, 1)$ ,  $\delta_{i,j\pm 1/2}(t) \in (0, 1)$ .*

*Proof.* The form of the flux comes from the computations present in this section. If we equate (17) to zero we can guarantee that the exact flux vanishes in the derived numerical approximation in the case where the components of the analytical flux vanish in the presence of a steady state. Finally, the latter result follows from the inequality  $\exp(x) \geq 1 + x$ .  $\square$

*Remark 2.* We can observe that for  $\lambda_{i+1/2,j}(t) \ll 1$  and  $\lambda_{i,j+1/2}(t) \ll 1$  we have

$$\delta_{i+1/2,j}(t) = \frac{1}{2} + O(\lambda_{i+1/2,j}(t)), \quad \delta_{i,j+1/2}(t) = \frac{1}{2} + O(\lambda_{i,j+1/2}(t)),$$

and, therefore, when  $\lambda_{i+1/2,j}(t) = \lambda_{i,j+1/2}(t) = 0$ , we have that  $\delta_{i+1/2,j}(t) = \delta_{i,j+1/2}(t) = \frac{1}{2}$ .

In the scheme defined by the numerical fluxes (21) with  $\tilde{\mathcal{G}}_{i+1/2,j}^x(t)$  and  $\tilde{\mathcal{G}}_{i,j+1/2}^y(t)$  defined in (22) and with  $\delta_{i+1/2,j}(t)$ ,  $\delta_{i,j+1/2}(t)$  defined by (25)-(26), we shall approximate the integrals appearing in the nonlinear weights (26) through high order quadrature rules, as done in [34]. In fact, it is worth observing that the derived scheme may be seen as a generalization of the classical second-order Chang-Cooper scheme [14, 25] to Fokker-Planck equations with general diffusion matrix. In their original formulation, these works focused on linear Fokker-Planck equations with diagonal diffusion matrix and a recent generalization to the nonlinear case has been proposed in [34]. We highlight how the scheme proposed in [34] and the present scheme are coherent to the original one by approximating the functions (26) through a midpoint quadrature rule as follows

$$\begin{aligned} \lambda_{i+1/2,j}^{\text{mid}}(t) &= \int_{w_{i,j}}^{w_{i+1,j}} \frac{1}{\mathcal{D}^1(w)} \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right) dw_x \\ &= \frac{\Delta w}{\mathcal{D}_{i+1/2,j}^1} \left( \mathcal{C}_{i+1/2,j}^x - \frac{\mathbb{D}_{i+1/2,j}^{1,2}}{\mathbb{D}_{i+1/2,j}^{2,2}} \mathcal{C}_{i+1/2,j}^y \right), \\ \lambda_{i,j+1/2}^{\text{mid}}(t) &= \int_{w_{i,j}}^{w_{i,j+1}} \frac{1}{\mathcal{D}^2(w)} \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right) dw_y \\ &= \frac{\Delta w}{\mathcal{D}_{i,j+1/2}^2} \left( \mathcal{C}_{i,j+1/2}^y - \frac{\mathbb{D}_{i,j+1/2}^{2,1}}{\mathbb{D}_{i,j+1/2}^{1,1}} \mathcal{C}_{i,j+1/2}^x \right), \end{aligned}$$



leading to the following weights

$$\begin{aligned}\delta_{i+1/2,j}^{\text{mid}}(t) &= \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w \left( \mathcal{C}_{i+1/2,j}^x(t) - \frac{\mathbb{D}_{i+1/2,j}^{1,2}}{\mathbb{D}_{i+1/2,j}^{2,2}} \mathcal{C}_{i+1/2,j}^y(t) \right)} + \frac{1}{1 - \exp(\lambda_{i+1/2,j}^{\text{mid}}(t))}, \\ \delta_{i,j+1/2}^{\text{mid}}(t) &= \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w \left( \mathcal{C}_{i,j+1/2}^y(t) - \frac{\mathbb{D}_{i,j+1/2}^{2,1}}{\mathbb{D}_{i,j+1/2}^{1,1}} \mathcal{C}_{i,j+1/2}^x(t) \right)} + \frac{1}{1 - \exp(\lambda_{i,j+1/2}^{\text{mid}}(t))}.\end{aligned}$$

Hence, in the case  $\mathbb{D}^{1,2}(w) = \mathbb{D}^{2,1}(w) = 0$  we recover the classical formulation. Furthermore, we observe that if  $\mathcal{B}[f](w, t)$  does not depend on  $f(w, t)$  and has components which are first order polynomials, the midpoint rule gives an exact evaluation of the integrals in (26). More generally, to extend the introduced approach as in [34] we may consider standard high order quadrature rules for the computation of the nonlinear weights (26), see e.g. [15].

*Remark 3.* In the case  $\mathcal{B}[f](w, t) = B(w)$  the quasi-stationary formulation (9) becomes stationary, because  $\mathcal{C}(w)$  does not depend on  $f(w, t)$  anymore. Once we know the stationary state  $f^\infty(w)$ , we can compute the weights  $\delta_{i+1/2,j}(t)$ ,  $\delta_{i,j+1/2}(t)$  exactly. In fact, we have

$$\begin{aligned}\frac{f_{i+1,j}^\infty}{f_{i,j}^\infty} &= \exp \left\{ - \int_{w_{i,j}}^{w_{i+1,j}} \frac{1}{\mathcal{D}^1(w)} \left( \mathcal{C}^x(w, t) - \frac{\mathbb{D}^{1,2}(w)}{\mathbb{D}^{2,2}(w)} \mathcal{C}^y(w, t) \right) dw_x \right\} \\ &= \exp \left\{ -\lambda_{i+1/2,j}^\infty \right\} \\ \frac{f_{i,j+1}^\infty}{f_{i,j}^\infty} &= \exp \left\{ - \int_{w_{i,j}}^{w_{i,j+1}} \frac{1}{\mathcal{D}^2(w)} \left( \mathcal{C}^y(w, t) - \frac{\mathbb{D}^{2,1}(w)}{\mathbb{D}^{1,1}(w)} \mathcal{C}^x(w, t) \right) dw_y \right\} \\ &= \exp \left\{ -\lambda_{i,j+1/2}^\infty \right\},\end{aligned}$$

that define the following weights

$$\begin{aligned}\delta_{i+1/2,j}^\infty &= \frac{1}{\log f_{i,j}^\infty - \log f_{i+1,j}^\infty} + \frac{f_{i+1,j}^\infty}{f_{i+1,j}^\infty - f_{i,j}^\infty}, \\ \delta_{i,j+1/2}^\infty &= \frac{1}{\log f_{i,j}^\infty - \log f_{i,j+1}^\infty} + \frac{f_{i,j+1}^\infty}{f_{i,j+1}^\infty - f_{i,j}^\infty}.\end{aligned}\tag{27}$$

Using classical methods, as it is done for example in [29] for the linear Fokker-Planck equation with diagonal diffusion matrix, we can prove that the proposed scheme is consistent if the problem is linear and the flux vanishes at the steady state. In particular, using the same arguments of [29], it is possible to see that the stationary state is approximated with an order equal to the order of the quadrature rule. This is cannot be proved for problems whose stationary state does not make the flux vanish (see Test2).

*Remark 4.* If we consider the limit case in which the diffusion tensor tends to be singular and the elements of  $\nabla \cdot \mathbb{D}$  tend to vanish, we obtain

$$\delta_{i+1/2,j}(t) = \begin{cases} 0, & \mathcal{B}_{i+1/2,j}(t) > 0, \\ 1, & \mathcal{B}_{i+1/2,j}(t) < 0, \end{cases} \quad \delta_{i,j+1/2}(t) = \begin{cases} 0 & \mathcal{B}_{i,j+1/2}(t) > 0, \\ 1 & \mathcal{B}_{i,j+1/2}(t) < 0. \end{cases}$$

Therefore the scheme reduces to a first order upwind scheme.

### 3 Main properties

In this section we show the properties of the derived numerical schemes. In particular, we will prove how the present method enforces conservation of mass, nonnegativity of the numerical solution and correctly dissipates the entropy.

### 3.1 Conservation of mass

We notice that the no-flux boundary conditions

$$\mathcal{F}(w, t) \cdot \mathbf{n}(w) = 0, \quad w \in \partial\Omega$$

ensure the conservation of mass in the problem (1) since

$$\int_{\Omega} f(w, t) dw = \int_{\Omega} f_0(w) dw$$

for all time  $t \geq 0$ . At the numerical level, no-flux boundary conditions are obtained by imposing

$$\mathcal{F}_{N+1/2, j}^x(t) = \mathcal{F}_{-1/2, j}^x(t) = 0, \quad \text{and} \quad \mathcal{F}_{i, N+1/2}^y(t) = \mathcal{F}_{i, -1/2}^y(t) = 0, \quad \forall i, j = 0, \dots, N, \quad t \geq 0, \quad (28)$$

and we can prove that the introduced scheme ensures the conservation of mass.

**Lemma 5.** *Let us consider a discretization of the problem (1) in the form (14) complemented with no-flux boundary conditions (28). Then we have*

$$\frac{d}{dt} \sum_{i=0}^N \sum_{j=0}^N f_{i,j}(t) = 0.$$

*Proof.* From (14) we have

$$\sum_{i=0}^N \sum_{j=0}^N \frac{d}{dt} f_{i,j}(t) = \frac{1}{\Delta w} \sum_{j=0}^N \left( \mathcal{F}_{-1/2, j}^x(t) - \mathcal{F}_{N+1/2, j}^x(t) \right) + \frac{1}{\Delta w} \sum_{i=0}^N \left( \mathcal{F}_{i, -1/2}^y(t) - \mathcal{F}_{i, N+1/2}^y(t) \right),$$

from which we conclude using (28).  $\square$

### 3.2 Positivity of the explicit scheme

In this section we will provide results for non-negativity of the scheme with explicit time integration. We introduce the time discretization  $t^n = n\Delta t$ ,  $n = 0, \dots, N_T$  with  $\Delta t = T/N_T$  being  $T$  the final time. We first consider the simple forward Euler method

$$f_{i,j}^{n+1} = f_{i,j}^n + \Delta t \frac{\mathcal{F}_{i+1/2, j}^{x,n} - \mathcal{F}_{i-1/2, j}^{x,n}}{\Delta w} + \Delta t \frac{\mathcal{F}_{i, j+1/2}^{y,n} - \mathcal{F}_{i, j-1/2}^{y,n}}{\Delta w},$$

where  $f_{i,j}^n = f_{i,j}(t^n)$  and  $\mathcal{F}_{i+1/2, j}^{x,n}, \mathcal{F}_{i, j+1/2}^{y,n}$  are the numerical fluxes at time  $t^n$ , i.e.  $\mathcal{F}_{i+1/2, j}^{x,n} = \mathcal{F}_{i+1/2, j}^x(t^n)$ , and  $\mathcal{F}_{i, j+1/2}^{y,n} = \mathcal{F}_{i, j+1/2}^y(t^n)$ .

We can prove the following result

**Theorem 6.** *Under the time step restriction*

$$\Delta t \leq \frac{\Delta w^2}{2[(G_x + G_y)\Delta w + (D^1 + D^2)]} \quad (29)$$

where

$$G_x = \max_{i,j,n} |\tilde{\mathcal{G}}_{i+1/2, j}^{x,n}|, \quad G_y = \max_{i,j,n} |\tilde{\mathcal{G}}_{i, j+1/2}^{y,n}|,$$

and

$$D^1 = \max_{i,j} \mathcal{D}_{i+1/2, j}^1, \quad D^2 = \max_{i,j} \mathcal{D}_{i, j+1/2}^2,$$

the explicit scheme preserves nonnegativity, i.e.  $f_{i,j}^{n+1} \geq 0$  if  $f_{i,j}^n \geq 0$ .

*Proof.* We will adopt the structure of the scheme introduced in Theorem 1. In details, the scheme reads

$$\begin{aligned}
f_{i,j}^{n+1} = & f_{i,j}^n + \frac{\Delta t}{\Delta w} \left[ \left( \tilde{\mathcal{G}}_{i+1/2,j}^{x,n} (1 - \delta_{i+1/2,j}^n) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \right) f_{i+1,j}^n \right. \\
& - \left( -\tilde{\mathcal{G}}_{i+1/2,j}^{x,n} \delta_{i+1/2,j}^n + \tilde{\mathcal{G}}_{i-1/2,j}^{x,n} (1 - \delta_{i-1/2,j}^n) + \frac{\mathcal{D}_{i+1/2,j}^1 + \mathcal{D}_{i-1/2,j}^1}{\Delta w} \right) f_{i,j}^n \\
& + \left. \left( -\tilde{\mathcal{G}}_{i-1/2,j}^{x,n} \delta_{i-1/2,j}^n + \frac{\mathcal{D}_{i-1/2,j}^1}{\Delta w} \right) f_{i-1,j}^n \right] + \frac{\Delta t}{\Delta w} \left[ \left( \tilde{\mathcal{G}}_{i,j+1/2}^{y,n} (1 - \delta_{i,j+1/2}^n) + \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \right) f_{i,j+1}^n \right. \\
& - \left( -\tilde{\mathcal{G}}_{i,j+1/2}^{y,n} \delta_{i,j+1/2}^n + \tilde{\mathcal{G}}_{i,j-1/2}^{y,n} (1 - \delta_{i,j-1/2}^n) + \frac{\mathcal{D}_{i,j+1/2}^2 + \mathcal{D}_{i,j-1/2}^2}{\Delta w} \right) f_{i,j}^n \\
& + \left. \left( -\tilde{\mathcal{G}}_{i,j-1/2}^{y,n} \delta_{i,j-1/2}^n + \frac{\mathcal{D}_{i,j-1/2}^2}{\Delta w} \right) f_{i,j-1}^n \right].
\end{aligned}$$

This is a sum of convex combinations of  $f_{i+1,j}^n$ ,  $f_{i-1,j}^n$  and  $f_{i,j+1}^n, f_{i,j-1}^n$  if the following conditions are satisfied

$$\begin{aligned}
\tilde{\mathcal{G}}_{i+1/2,j}^{x,n} (1 - \delta_{i+1/2,j}^n) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} &\geq 0, & -\tilde{\mathcal{G}}_{i-1/2,j}^{x,n} \delta_{i-1/2,j}^n + \frac{\mathcal{D}_{i-1/2,j}^1}{\Delta w} &\geq 0, \\
\tilde{\mathcal{G}}_{i,j+1/2}^{y,n} (1 - \delta_{i,j+1/2}^n) + \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} &\geq 0, & -\tilde{\mathcal{G}}_{i,j-1/2}^{y,n} \delta_{i,j-1/2}^n + \frac{\mathcal{D}_{i,j-1/2}^2}{\Delta w} &\geq 0,
\end{aligned}$$

that is equivalent to

$$\begin{aligned}
\lambda_{i+1/2,j}^n \left( 1 - \frac{1}{1 - \exp(\lambda_{i+1/2,j}^n)} \right) &\geq 0, & \frac{\lambda_{i-1/2,j}^n}{\exp(\lambda_{i-1/2,j}^n) - 1} &\geq 0, \\
\lambda_{i,j+1/2}^n \left( 1 - \frac{1}{1 - \exp(\lambda_{i,j+1/2}^n)} \right) &\geq 0, & \frac{\lambda_{i,j-1/2}^n}{\exp(\lambda_{i,j-1/2}^n) - 1} &\geq 0,
\end{aligned}$$

which hold true thanks to the basic properties of the exponential function. In order to ensure positivity for  $f_{i,j}^{n+1}$  if  $f_{i,j}^n \geq 0$  we must have for all  $i, j$

$$\left( 1 - (\nu_x + \nu_y) \frac{\Delta t}{\Delta w} \right) f_{i,j}^n \geq 0,$$

where

$$\begin{aligned}
\nu_x = \max_{i,j} \left\{ -\tilde{\mathcal{G}}_{i+1/2,j}^{x,n} \delta_{i+1/2,j}^n + \tilde{\mathcal{G}}_{i-1/2,j}^{x,n} (1 - \delta_{i-1/2,j}^n) + \frac{\mathcal{D}_{i+1/2,j}^1 + \mathcal{D}_{i-1/2,j}^1}{\Delta w} \right\}, \\
\nu_y = \max_{i,j} \left\{ -\tilde{\mathcal{G}}_{i,j+1/2}^{y,n} \delta_{i,j+1/2}^n + \tilde{\mathcal{G}}_{i,j-1/2}^{y,n} (1 - \delta_{i,j-1/2}^n) + \frac{\mathcal{D}_{i,j+1/2}^2 + \mathcal{D}_{i,j-1/2}^2}{\Delta w} \right\},
\end{aligned}$$

from which we can conclude as  $0 \leq \delta_{i\pm 1/2,j} \leq 1$ ,  $0 \leq \delta_{i,j\pm 1/2} \leq 1$ .  $\square$

We highlight how the CFL criterion in (29) ensures positivity of the numerical solution of the problem. This remarkable property holds also for higher order strong stability preserving (SSP) methods like SSP Runge-Kutta and SSP multistep methods [23] since these are convex combinations of the forward Euler integration. Hence, the proved non-negativity of the scheme is automatically extended to each general SSP type time integration.

Even if in (29) we obtained an effective time step bound for the positivity of the explicit scheme, for practical purposes this parabolic restriction is very heavy especially in genuine nonlinear type problems. Usually the strategy to overcome this problem relies in the adoption of IMEX schemes [17]. Nevertheless, this is not always possible due to the strong nonlinearities embedded in problem (1) coming from the nonlocal drift term. **Furthermore**, the defined weights  $\delta_{i+1/2,j}$ ,  $\delta_{i,j+1/2}$  depend in general on  $f$  introducing additional difficulties. An efficient way to overcome this problem relies in the semi-implicit integration technique, see [7].

### 3.3 Positivity of the semi-implicit scheme

To apply semi-implicit schemes we integrate (14) as follows

$$f_{i,j}^{n+1} = f_{i,j}^n + \Delta t \frac{\hat{\mathcal{F}}_{i+1/2,j}^{x,n+1} - \hat{\mathcal{F}}_{i-1/2,j}^{x,n+1}}{\Delta w} + \Delta t \frac{\hat{\mathcal{F}}_{i,j+1/2}^{y,n+1} - \hat{\mathcal{F}}_{i,j-1/2}^{y,n+1}}{\Delta w}, \quad (30)$$

where now the discretized flux terms  $\hat{\mathcal{F}}_{i+1/2,j}^{x,n+1}$ ,  $\hat{\mathcal{F}}_{i,j+1/2}^{y,n+1}$  are defined as

$$\begin{aligned} \hat{\mathcal{F}}_{i+1/2,j}^{x,n+1} &= \tilde{\mathcal{G}}_{i+1/2,j}^{x,n} \left[ (1 - \delta_{i+1/2,j}^n) f_{i+1,j}^{n+1} + \delta_{i+1/2,j}^n f_{i,j}^{n+1} \right] + \mathcal{D}_{i+1/2,j}^1 \frac{f_{i+1,j}^{n+1} - f_{i,j}^{n+1}}{\Delta w}, \\ \hat{\mathcal{F}}_{i,j+1/2}^{y,n+1} &= \tilde{\mathcal{G}}_{i,j+1/2}^{y,n} \left[ (1 - \delta_{i,j+1/2}^n) f_{i,j+1}^{n+1} + \delta_{i,j+1/2}^n f_{i,j}^{n+1} \right] + \mathcal{D}_{i,j+1/2}^2 \frac{f_{i,j+1}^{n+1} - f_{i,j}^{n+1}}{\Delta w}. \end{aligned}$$

**Theorem 7.** *Under the time step restriction*

$$\Delta t \leq \frac{\Delta w}{2(G_x + G_y)}, \quad G_x = \max_{i,j,n} \{|\tilde{\mathcal{G}}_{i+1/2,j}^{x,n}|\}, \quad G_y = \max_{i,j,n} \{|\tilde{\mathcal{G}}_{i,j+1/2}^{y,n}|\},$$

*the semi-implicit scheme (30) preserves nonnegativity, i.e.,  $f_{i,j}^{n+1} \geq 0$  if  $f_{i,j}^n \geq 0$ .*

*Proof.* Equation (30) corresponds to

$$\begin{aligned} & f_{i+1,j}^{n+1} \left\{ -\frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i+1/2,j}^{x,n} (1 - \delta_{i+1/2,j}^n) + \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \right] \right\} \\ & + f_{i,j}^{n+1} \left\{ 1 - \frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i+1/2,j}^{x,n} \delta_{i+1/2,j}^n - \tilde{\mathcal{G}}_{i-1/2,j}^{x,n} (1 - \delta_{i-1/2,j}^n) - \frac{\mathcal{D}_{i+1/2,j}^1 + \mathcal{D}_{i-1/2,j}^1}{\Delta w} \right] \right\} \\ & + f_{i-1,j}^{n+1} \left\{ -\frac{\Delta t}{\Delta w} \left[ -\tilde{\mathcal{G}}_{i-1/2,j}^{x,n} \delta_{i-1/2,j}^n + \frac{\mathcal{D}_{i-1/2,j}^1}{\Delta w} \right] \right\} \\ & + f_{i,j+1}^{n+1} \left\{ -\frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i,j+1/2}^{y,n} (1 - \delta_{i,j+1/2}^n) + \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \right] \right\} \\ & + f_{i,j}^{n+1} \left\{ 1 - \frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i,j+1/2}^{y,n} \delta_{i,j+1/2}^n - \tilde{\mathcal{G}}_{i,j-1/2}^{y,n} (1 - \delta_{i,j-1/2}^n) - \frac{\mathcal{D}_{i,j+1/2}^2 + \mathcal{D}_{i,j-1/2}^2}{\Delta w} \right] \right\} \\ & + f_{i,j-1}^{n+1} \left\{ -\frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i,j-1/2}^{y,n} \delta_{i,j-1/2}^n + \frac{\mathcal{D}_{i,j-1/2}^2}{\Delta w} \right] \right\} = f_{i,j}^n. \end{aligned}$$

Using the identities in (26), we have that

$$\begin{aligned}
& f_{i+1,j}^{n+1} \left\{ -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i+1/2,j}^1 \frac{\lambda_{i+1/2,j}^n}{\exp(\lambda_{i+1/2,j}^n) - 1} \exp(\lambda_{i+1/2,j}^n) \right\} \\
& + f_{i,j}^{n+1} \left\{ 1 + \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i+1/2,j}^1 \frac{\lambda_{i+1/2,j}^n}{\exp(\lambda_{i+1/2,j}^n) - 1} + \mathcal{D}_{i-1/2,j}^1 \frac{\lambda_{i-1/2,j}^n}{\exp(\lambda_{i-1/2,j}^n) - 1} \exp(\lambda_{i-1/2,j}^n) \right] \right\} \\
& + f_{i-1,j}^{n+1} \left\{ -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i-1/2,j}^1 \frac{\lambda_{i-1/2,j}^n}{\exp(\lambda_{i-1/2,j}^n) - 1} \right\} \\
& + f_{i,j+1}^{n+1} \left\{ -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i,j+1/2}^2 \frac{\lambda_{i,j+1/2}^n}{\exp(\lambda_{i,j+1/2}^n) - 1} \exp(\lambda_{i,j+1/2}^n) \right\} \\
& + f_{i,j}^{n+1} \left\{ 1 + \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i,j+1/2}^2 \frac{\lambda_{i,j+1/2}^n}{\exp(\lambda_{i,j+1/2}^n) - 1} + \mathcal{D}_{i,j-1/2}^2 \frac{\lambda_{i,j-1/2}^n}{\exp(\lambda_{i,j-1/2}^n) - 1} \exp(\lambda_{i,j-1/2}^n) \right] \right\} \\
& + f_{i,j-1}^{n+1} \left\{ -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i,j-1/2}^2 \frac{\lambda_{i,j-1/2}^n}{\exp(\lambda_{i,j-1/2}^n) - 1} \right\} = f_{i,j}^n.
\end{aligned}$$

Then by introducing the quantities

$$\alpha_{i+1/2,j}^n = \frac{\lambda_{i+1/2,j}^n}{\exp(\lambda_{i+1/2,j}^n) - 1} \geq 0 \quad \text{and} \quad \alpha_{i,j+1/2}^n = \frac{\lambda_{i,j+1/2}^n}{\exp(\lambda_{i,j+1/2}^n) - 1} \geq 0$$

and setting

$$\begin{aligned}
R_x(j)_i^n &= 1 + \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i+1/2,j}^1 \alpha_{i+1/2,j}^n - \mathcal{D}_{i-1/2,j}^1 \alpha_{i-1/2,j}^n \exp(\lambda_{i-1/2,j}^n) \right], \\
Q_x(j)_i^n &= -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i+1/2,j}^1 \alpha_{i+1/2,j}^n \exp(\lambda_{i+1/2,j}^n), \\
P_x(j)_i^n &= -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i-1/2,j}^1 \alpha_{i-1/2,j}^n, \\
R_y(i)_j^n &= 1 + \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i,j+1/2}^2 \alpha_{i,j+1/2}^n - \mathcal{D}_{i,j-1/2}^2 \alpha_{i,j-1/2}^n \exp(\lambda_{i,j-1/2}^n) \right], \\
Q_y(i)_j^n &= -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i,j+1/2}^2 \alpha_{i,j+1/2}^n \exp(\lambda_{i,j+1/2}^n), \\
P_y(i)_j^n &= -\frac{\Delta t}{\Delta w^2} \mathcal{D}_{i,j-1/2}^2 \alpha_{i,j-1/2}^n,
\end{aligned}$$

the latter equation reduces to

$$\begin{aligned}
& R_x(j)_i^n f_{i,j}^{n+1} - Q_x(j)_i^n f_{i+1,j}^{n+1} - P_x(j)_i^n f_{i-1,j}^{n+1} \\
& + R_y(i)_j^n f_{i,j}^{n+1} - Q_y(i)_j^n f_{i,j+1}^{n+1} - P_y(i)_j^n f_{i,j-1}^{n+1} = f_{i,j}^n.
\end{aligned}$$

Now, by denoting  $\mathbf{f}^n = \{f_{i,j}^n\}_{i=1,\dots,N}^{j=1,\dots,N}$  we can define the matrices

$$\begin{aligned}
\mathcal{A}_x[\mathbf{f}^n]_{ik} &= \begin{cases} R_x(j)_i^n & k = i, \\ -Q_x(j)_i^n & k = i + 1, \quad 0 \leq i \leq N - 1, \\ -P_x(j)_i^n & k = i - 1, \quad 1 \leq i \leq N, \end{cases} \\
\mathcal{A}_y[\mathbf{f}^n]_{jk} &= \begin{cases} R_y(i)_j^n & k = j, \\ -Q_y(i)_j^n & k = j + 1, \quad 0 \leq j \leq N - 1, \\ -P_y(i)_j^n & k = j - 1, \quad 1 \leq j \leq N, \end{cases}
\end{aligned}$$

and we reduce to study

$$(\mathcal{A}_x[\mathbf{f}^n] + \mathcal{A}_y[\mathbf{f}^n]) \mathbf{f}^{n+1} = \mathbf{f}^n.$$

If  $\mathbf{f}^n \geq 0$ , in order to prove that  $\mathbf{f}^{n+1} \geq 0$  it is sufficient to prove that  $(\mathcal{A}_x[\mathbf{f}^n] + \mathcal{A}_y[\mathbf{f}^n])^{-1}$  is nonnegative. Let us observe that since  $(\mathcal{A}_x[\mathbf{f}^n] + \mathcal{A}_y[\mathbf{f}^n])$  is tridiagonal we only need to prove that it is a diagonally dominant matrix. In particular, this is true if for each  $i, j = 1, \dots, N$  the following inequality is verified

$$|R_x(j)_i^n + R_y(i)_j^n| > |Q_x(j)_i^n + Q_y(i)_j^n| + |P_x(j)_i^n + P_y(i)_j^n|,$$

which is true provided

$$\begin{aligned} 1 &> \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i+1/2,j}^1 \alpha_{i+1/2,j}^n (\exp(\lambda_{i+1/2,j}^n) - 1) - \mathcal{D}_{i-1/2,j}^1 \alpha_{i-1/2,j}^n (\exp(\lambda_{i-1/2,j}^n) - 1) \right] \\ &\quad + \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i,j+1/2}^2 \alpha_{i,j+1/2}^n (\exp(\lambda_{i,j+1/2}^n) - 1) - \mathcal{D}_{i,j-1/2}^2 \alpha_{i,j-1/2}^n (\exp(\lambda_{i,j-1/2}^n) - 1) \right] \\ &= \frac{\Delta t}{\Delta w^2} \left[ \mathcal{D}_{i+1/2,j}^1 \lambda_{i+1/2,j}^n - \mathcal{D}_{i-1/2,j}^1 \lambda_{i-1/2,j}^n + \mathcal{D}_{i,j+1/2}^2 \lambda_{i,j+1/2}^n - \mathcal{D}_{i,j-1/2}^2 \lambda_{i,j-1/2}^n \right] \\ &= \frac{\Delta t}{\Delta w} \left[ \tilde{\mathcal{G}}_{i+1/2,j}^{x,n} - \tilde{\mathcal{G}}_{i-1/2,j}^{x,n} + \tilde{\mathcal{G}}_{i,j+1/2}^{y,n} - \tilde{\mathcal{G}}_{i,j-1/2}^{y,n} \right]. \end{aligned}$$

□

*Remark 8.* Fully-implicit schemes require a special treatment since the nonlinearity in the drift term poses nontrivial questions at the numerical level. A possible way to overcome this difficulty is to use iterative methods as suggested in [34]. This issue anyway goes beyond the goals of the present manuscript and we postpone this discussion to future works.

## 4 Trend to equilibrium

A classical question in kinetic theory pertains to the determination of the rate of exponential convergence to equilibrium. To this end the consolidated approach relies on entropy production arguments for which lower bounds are explicitly computable thanks to log-Sobolev inequalities, see [40, 42]. In particular, the convergence to the stationary state of the standard Fokker-Planck equation can be achieved by looking at the monotonicity in time of various Lyapunov functionals like the relative entropy. In the nonconstant diffusion case additional difficulties arise since standard log-Sobolev inequality are not available [28].

### 4.1 Steady state and vanishing flux for linear problems

In order to study the entropy properties, as done typically [40, 42, 21], we consider the linear problem defined by  $\mathcal{B}[f](w, t) = B(w)$ . Moreover, we suppose that a stationary state exists and that, coherently with the assumptions already discussed, the flux vanishes at the stationary state, i.e.  $\mathcal{F}^\infty(w) = 0$ . The latter is equivalent to say that  $f^\infty(w)$  satisfies

$$B(w) f^\infty(w) + \nabla_w \cdot (\mathbb{D}(w) f^\infty(w)) = 0, \quad w \in \Omega.$$

Then we have

$$B(w) = - \frac{f^\infty(w) \nabla_w \cdot \mathbb{D}(w)}{f^\infty(w)} - \mathbb{D}(w) \frac{\nabla_w f^\infty(w)}{f^\infty(w)} = - \nabla_w \cdot \mathbb{D}(w) - \mathbb{D}(w) \frac{\nabla_w f^\infty(w)}{f^\infty(w)}, \quad (31)$$

see [37]. Hence, we can rewrite our problem in the form

$$\partial_t f(w, t) = \nabla_w \cdot \left[ f^\infty(w) \mathbb{D}(w) \nabla_w \frac{f(w, t)}{f^\infty(w)} \right], \quad (32)$$

since

$$\begin{aligned}
& \nabla_w \cdot [B(w)f(w, t) + \nabla_w \cdot (\mathbb{D}(w)f(w, t))] \\
&= \nabla_w \cdot \left[ -f(w, t) \nabla_w \cdot \mathbb{D}(w) - f(w, t) \mathbb{D}(w) \frac{\nabla_w f^\infty(w)}{f^\infty(w)} + \nabla_w \cdot (\mathbb{D}(w)f(w, t)) \right] \\
&= \nabla_w \cdot \left[ -f(w, t) \mathbb{D}(w) \frac{\nabla_w f^\infty(w)}{f^\infty(w)} + \mathbb{D}(w) \nabla_w f(w, t) \right] \\
&= \nabla_w \cdot \left[ f(w, t) \mathbb{D}(w) \left( \frac{\nabla_w f(w, t)}{f(w, t)} - \frac{\nabla_w f^\infty(w)}{f^\infty(w)} \right) \right] \\
&= \nabla_w \cdot \left[ f(w, t) \mathbb{D}(w) \nabla_w \log \left( \frac{f(w, t)}{f^\infty(w)} \right) \right] \\
&= \nabla_w \cdot \left[ f^\infty(w) \mathbb{D}(w) \nabla_w \frac{f(w, t)}{f^\infty(w)} \right].
\end{aligned}$$

The no-flux boundary conditions in this case read

$$\left[ f^\infty(w) \mathbb{D}(w) \nabla_w \frac{f(w, t)}{f^\infty(w)} \right] \cdot \mathbf{n}(w) = 0, \quad w \in \partial\Omega.$$

Therefore, from the Landau's formulation (32), we get the equation satisfied by  $F(w, t) = f(w, t)/f^\infty(w)$  that is

$$\begin{aligned}
\partial_t F(w, t) &= \frac{\partial_t f(w, t)}{f^\infty(w)} = \frac{\nabla_w \cdot [f^\infty(w) \mathbb{D}(w) \nabla_w F(w, t)]}{f^\infty(w)} \\
&= \nabla_w \cdot (\mathbb{D}(w) \nabla_w F(w, t)) + (\mathbb{D}(w) \nabla_w F) \cdot \frac{\nabla_w f^\infty(w)}{f^\infty(w)} \\
&= \nabla_w \cdot (\mathbb{D}(w) \nabla_w F(w, t)) - B(w) \cdot \nabla_w F(w, t) - (\nabla_w \cdot \mathbb{D}(w)) \cdot \nabla_w F(w, t),
\end{aligned}$$

where the last equality holds true since  $\mathbb{D}(w)$  is a symmetric matrix  $\forall w \in \Omega$  and thanks to the relation (31). Now, since

$$\nabla_w \cdot (\mathbb{D}(w) \nabla_w F(w, t)) = (\nabla_w \cdot \mathbb{D}(w)) \cdot \nabla_w F(w, t) + \mathbb{D}(w) : \nabla_w (\nabla_w F(w, t)), \quad (33)$$

where  $\nabla_w (\nabla_w F(w, t))$  is the covariant derivative of the vector  $\nabla_w F(w, t)$ , i.e.  $\nabla_w (\nabla_w F(w, t)) = (\partial_{w_i} \nabla_w F(w, t)) = (\partial_{w_i} \partial_{w_j} F(w, t))$ , and it is the Hessian matrix of  $F$ , which we will denote  $H_w[F]$ . With : we indicated the inner tensorial product that is for definition

$$\mathbb{D}(w) : H_w[F](w, t) = \text{tr} \left[ (H_w[F](w, t))^T \mathbb{D}(w) \right].$$

In conclusion, we obtain

$$\partial_t F(w, t) = \mathbb{D}(w) : H_w[F](w, t) - B(w) \cdot \nabla_w F(w, t). \quad (34)$$

## 4.2 Lyapunov functionals

We will focus on the study of relative Shannon entropy for the problem (1) with nonconstant diffusion. We will extend the results proved in [21] to the two-dimensional case where the diffusion is a nonconstant positive definite tensor of the second order and the drift term is general in the form  $B(w)$ .

Let  $f, g : \Omega \mapsto \mathbb{R}^+$  denote two probability densities. Then, the relative Shannon entropy of  $f$  and  $g$  is defined by

$$H(f|g) = \int_{\Omega} f \log \frac{f}{g} dw. \quad (35)$$

It is a Lyapunov functional since the following result can be established.

**Theorem 9.** Let us suppose that  $\mathcal{F}^\infty(w) = 0$  and that the drift is of the form (31). Let  $F(w, t)$  be the solution of (34) in  $\Omega$ . Then, if  $\Psi(w)$  is a smooth function such that

$$|\Psi| \leq c \leq \infty \quad \text{on } \partial\Omega,$$

then the following relation holds

$$\int_{\Omega} f^\infty(w, t) \Psi(w) \partial_t F(w, t) dw = \int_{\Omega} f^\infty(w, t) \nabla_w \Psi \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw.$$

*Proof.* From (34) it follows that

$$\int_{\Omega} f^\infty(w) \Psi(w) \partial_t F(w, t) dw = \int_{\Omega} f^\infty(w) \Psi(w) (\mathbb{D}(w) : H_w[F](w, t) - B(w) \cdot \nabla_w F(w, t)) dw$$

and from (35) the latter term is equal to

$$\begin{aligned} & \int_{\Omega} f^\infty(w) \Psi(w) \left[ \nabla_w (\mathbb{D}(w) \nabla_w F(w, t)) - \nabla_w \cdot \mathbb{D}(w) \nabla_w F(w, t) \right] dw - \int_{\Omega} f^\infty(w) \Psi(w) B(w) \cdot \nabla_w F(w, t) dw \\ &= - \int_{\Omega} \nabla_w (f^\infty(w) \Psi(w)) \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw + \oint_{\partial\Omega} \Psi(w) f^\infty(w) (\mathbb{D}(w) \nabla_w F(w, t)) \cdot \mathbf{n}(w) d\sigma(w) \\ &\quad - \int_{\Omega} \left[ B(w) f^\infty(w) + \nabla_w \cdot \mathbb{D}(w) f^\infty(w) \right] \cdot \nabla_w F(w, t) \Psi(w) dw \\ &= - \int_{\Omega} \nabla_w (f^\infty(w) \Psi(w)) \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw \\ &\quad - \int_{\Omega} \left[ B(w) f^\infty(w) + \nabla_w \cdot \mathbb{D}(w) f^\infty(w) \right] \cdot \nabla_w F(w, t) \Psi(w) dw \\ &= - \int_{\Omega} \Psi(w) \nabla_w f^\infty(w) \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw - \int_{\Omega} f^\infty(w) \nabla_w \Psi(w) \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw \\ &\quad - \int_{\Omega} \left[ B(w) f^\infty(w) + \nabla_w \cdot \mathbb{D}(w) f^\infty(w) \right] \cdot \nabla_w F(w, t) \Psi(w) dw \\ &= - \int_{\Omega} f^\infty(w) \nabla_w \Psi(w) \cdot (\mathbb{D}(w) \nabla_w F(w, t)) dw \\ &\quad - \int_{\Omega} \left[ B(w) f^\infty(w) + \nabla_w \cdot (\mathbb{D}(w, t) f^\infty(w)) \right] \cdot \nabla_w F(w, t) \Psi(w) dw \\ &= - \int_{\Omega} f^\infty(w) \nabla_w \Psi(w) \cdot (\mathbb{D}(w, t) \nabla_w F(w, t)) dw, \end{aligned}$$

as the border terms vanish because of the boundary conditions and where we used (33), the divergence theorem and the fact that  $\mathcal{F}^\infty(w) = 0$ .  $\square$

**Theorem 10.** Let us suppose that  $\mathcal{F}^\infty(w) = 0$  and that the drift is of the form (31). Let the smooth function  $\Phi(x)$ ,  $x \in \mathbb{R}^+$  be convex. Then, if  $F(w, t)$  is the solution of (34) in  $\Omega$ , and  $c \leq F(w, t) \leq C$  for some positive constants  $c < C$ , the functional

$$\Theta[F](t) = \int_{\Omega} f^\infty(w) \Phi(F(w, t)) dw$$

is monotonically decreasing in time, and the following equality holds

$$\frac{d}{dt} \Theta[F](t) = -I_\Theta[F](t),$$

where  $I_\Theta[F](t)$  denotes the quantity

$$I_\Theta[F](t) = \int_{\Omega} f^\infty(w) \Phi''(F(w, t)) \nabla_w F(w, t) \mathbb{D}(w) \nabla_w F(w, t) dw, \quad (36)$$

that is non-negative because  $\Phi$  is convex and  $\mathbb{D}(w)$  is positive definite.



*Proof.* The relation (36) follows from Theorem 9 by choosing  $\Psi(w) = \Phi'(F(w, t))$  for a fixed  $t > 0$ .  $\square$

The Shannon entropy of  $f(w, t)$  relative to  $f^\infty(w)$ , defined by (35) with  $g = f^\infty$ , is obtained by choosing  $\Phi(x) = x \log x$ . In this case

$$I_\Theta[F](t) = \int_\Omega f(w, t) \frac{\nabla_w F(w, t)}{F(w, t)} \mathbb{D}(w) \frac{\nabla_w F(w, t)}{F(w, t)} dw,$$

that may be re-written as

$$I_\Theta[F](t) = \int_\Omega f(w, t) \left( \frac{\nabla_w f(w, t)}{f(w, t)} - \frac{\nabla_w f^\infty(w)}{f^\infty(w)} \right) \mathbb{D}(w) \left( \frac{\nabla_w f(w, t)}{f(w, t)} - \frac{\nabla_w f^\infty(w)}{f^\infty(w)} \right) dw,$$

that is the Fisher information of  $f(w, t)$  relative to  $f^\infty(w)$ . We might also consider the weighted  $L^2$  distance that is obtained by considering  $\Phi(x) = (x - 1)^2$ . In this case

$$\Theta[F](t) = L^2(f, f^\infty) = \int_\Omega \frac{(f(w, t) - f^\infty(w))^2}{f^\infty(w)} dw$$

and

$$I(\Theta)[F](t) = 2 \int_\Omega \nabla_w F(w, t) \mathbb{D}(w) \nabla_w F(w, t) dw.$$

### 4.3 Dissipation of the numerical entropy

In the following results we show how the derived schemes dissipate in the introduced setting a Shannon-type numerical entropy functional.

**Theorem 11.** *Let us consider a drift term of the form (31). The numerical flux function (21) with  $\delta_{i+1/2,j}, \delta_{i,j+1/2}$  given by (25) can be written in the form (32) and reads*

$$\begin{cases} \mathcal{F}_{i+1/2,j}^x(t) = \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \hat{f}_{i+1/2,j}^\infty \left( \frac{f_{i+1,j}(t)}{f_{i+1,j}^\infty} - \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right), \\ \mathcal{F}_{i,j+1/2}^y(t) = \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \hat{f}_{i,j+1/2}^\infty \left( \frac{f_{i,j+1}(t)}{f_{i,j+1}^\infty} - \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right), \end{cases}$$

where

$$\hat{f}_{i+1/2,j}^\infty = \frac{f_{i+1,j}^\infty f_{i,j}^\infty}{f_{i+1,j}^\infty - f_{i,j}^\infty} \log \left( \frac{f_{i+1,j}^\infty}{f_{i,j}^\infty} \right), \quad \hat{f}_{i,j+1/2}^\infty = \frac{f_{i,j+1}^\infty f_{i,j}^\infty}{f_{i,j+1}^\infty - f_{i,j}^\infty} \log \left( \frac{f_{i,j+1}^\infty}{f_{i,j}^\infty} \right).$$

*Proof.* If  $\mathcal{B} = B(w)$ , we have that the definitions of  $\lambda_{i+1/2,j}$  and  $\lambda_{i,j+1/2}$  do not depend on time. Hence, we may denote  $\lambda_{i+1/2,j} = \lambda_{i+1/2,j}^\infty$  and  $\lambda_{i,j+1/2} = \lambda_{i,j+1/2}^\infty$  and we have

$$\log f_{i+1,j}^\infty - \log f_{i,j}^\infty = \lambda_{i+1/2,j},$$

$$\log f_{i,j+1}^\infty - \log f_{i,j}^\infty = \lambda_{i,j+1/2},$$

and  $\delta_{i+1/2,j}, \delta_{i,j+1/2}$  are of the form (27). Therefore, under these assumptions the flux function writes

$$\begin{aligned} \mathcal{F}_{i+1/2,j}^x(t) &= \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \left( \lambda_{i+1/2,j} \tilde{f}_{i+1/2,j}(t) + (f_{i+1,j}(t) - f_{i,j}(t)) \right) \\ &= \frac{\mathcal{D}_{i+1/2,j}^1}{\Delta w} \left( \lambda_{i+1/2,j} (f_{i+1,j}(t) + \delta_{i+1/2,j}(t) (f_{i,j}(t) - f_{i+1,j}(t))) + (f_{i+1,j}(t) - f_{i,j}(t)) \right) \end{aligned} \tag{37}$$

and

$$\begin{aligned}\mathcal{F}_{i,j+1/2}^y(t) &= \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \left( \lambda_{i,j+1/2} \tilde{f}_{i,j+1/2}(t) + (f_{i,j+1}(t) - f_{i,j}(t)) \right) \\ &= \frac{\mathcal{D}_{i,j+1/2}^2}{\Delta w} \left( \lambda_{i,j+1/2} (f_{i,j+1}(t) + \delta_{i,j+1/2} (f_{i,j}(t) - f_{i,j+1}(t))) + (f_{i,j+1}(t) - f_{i,j}(t)) \right).\end{aligned}\tag{38}$$

By substituting (27) in (37)-(38) we obtain the thesis.  $\square$

**Theorem 12.** *Let us consider a drift term of the form (31). The numerical flux (21) satisfies the discrete entropy dissipation*

$$\frac{d}{dt} \mathcal{H}_\Delta(f, f^\infty)(t) = -\mathcal{I}_\Delta(f, f^\infty)(t),$$

where

$$\mathcal{H}_\Delta(f, f^\infty)(t) = \Delta w^2 \sum_{j=0}^N \sum_{i=0}^N f_{i,j}(t) \log \frac{f_{i,j}(t)}{f_{i,j}^\infty}$$

and  $\mathcal{I}_\Delta$  is the positive discrete dissipation function

$$\begin{aligned}\mathcal{I}_\Delta(t) &= \Delta w \sum_{j=0}^N \sum_{i=0}^N \left[ \log \left( \frac{f_{i+1,j}(t)}{f_{i+1,j}^\infty} \right) - \log \left( \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \right] \left( \frac{f_{i+1,j}(t)}{f_{i+1,j}^\infty} - \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \hat{f}_{i+1/2,j}^\infty \mathcal{D}_{i+1/2,j}^1 \\ &\quad + \sum_{i=0}^N \sum_{j=0}^N f_{i,j+1}(t) \left[ \log \left( \frac{f_{i,j+1}(t)}{f_{i,j+1}^\infty} \right) - \log \left( \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \right] \left( \frac{f_{i,j+1}(t)}{f_{i,j+1}^\infty} - \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \hat{f}_{i,j+1/2}^\infty \mathcal{D}_{i,j+1/2}^2.\end{aligned}\tag{39}$$

*Proof.* If we compute the time derivative of the discrete relative entropy we have that

$$\begin{aligned}\frac{d}{dt} \mathcal{H}_\Delta(f, f^\infty)(t) &= \Delta w^2 \sum_{j=0}^N \sum_{i=0}^N \frac{df_{i,j}(t)}{dt} \left( 1 + \log \left( \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \right) \\ &= \Delta w \sum_{j=0}^N \sum_{i=0}^N \left( 1 + \log \left( \frac{f_{i,j}(t)}{f_{i,j}^\infty} \right) \right) \\ &\quad \times \left( \mathcal{F}_{i+1/2,j}^x(t) - \mathcal{F}_{i-1/2,j}^x(t) + \mathcal{F}_{i,j+1/2}^y(t) - \mathcal{F}_{i,j-1/2}^y(t) \right).\end{aligned}$$

After telescopic summation and thanks to the identity of Proposition 11 we obtain (39), which is positive because  $\mathcal{D}^\alpha > 0$ ,  $\alpha = 1, 2$  and  $(x - y) \log(\frac{x}{y})$  is positive for all  $x, y \geq 0$ .  $\square$

*Remark 13.* We highlight that in the case in which  $\mathbb{D}_{1,2}(w) = \mathbb{D}_{2,1}(w) = 0$  and  $\mathbb{D}(w)$  is isotropic, if we define an energy in the form

$$\xi(w, t) = (U_p * f)(w, t) + \frac{\text{tr}(\mathbb{D}(w))}{2} \log(f(w, t))$$

which in our case corresponds to

$$\mathcal{B}[f](w, t) = \nabla_w (U_p * f)(w, t),$$

with  $U_p = U_p(|w|)$  an interaction potential, then we have that

$$\nabla_w \xi(w, t) = \mathcal{B}[f](w, t) + \mathbb{D}(w) \nabla_w \log(f(w, t)).$$

Therefore, Eq. (1) may be written in the form

$$\partial_t f(w, t) = \nabla_w \cdot [f(w, t) \nabla_w \xi(w, t)], \quad w \in \Omega,$$

and therefore in a gradient flow structure for which entropic averaged schemes may be used [34].

|      | $SP_k$ |        |        |        | $SP_k$ |        |        |        |
|------|--------|--------|--------|--------|--------|--------|--------|--------|
| Time | 2      | 4      | 6      | G      | 2      | 4      | 6      | G      |
| 1    | 1.9601 | 1.6775 | 2.1106 | 2.111  | 1.9606 | 1.8176 | 2.1015 | 2.2103 |
| 10   | 1.9662 | 3.9708 | 7.4700 | 8.1449 | 1.9662 | 3.9708 | 7.4753 | 8.1449 |
| 20   | 1.9662 | 3.9708 | 7.4768 | 8.1453 | 1.9662 | 3.9708 | 7.4760 | 8.1449 |

Table 1: **Test 1.** Estimation of the order of convergence for  $SP_k$  scheme with explicit Euler (left) and RK4 (right). Rates have been computed using  $N = 21, 41, 81$  grid points in each component of the computational cell. We considered  $\sigma_1^2 = \sigma_2^2 = 1$ ,  $\rho = 0.1$ ,  $\Delta t = \Delta w^2 / (10\sigma_1^2 \Delta w + 10)$ .

## 5 Numerical tests

In this section we present some numerical examples of the class of Fokker-Planck [equations under study](#) with nonconstant full diffusion matrix solved through structure-preserving schemes that have been introduced in the previous sections. We will approximate the long time behaviour of (1) with  $d = 2$ , using the scheme defined by (21)-(22)-(25)-(26) with no-flux boundary conditions (28). [In the following, we will show numerically how the high order approximation of the nonlinear weights \(26\) reflects in an improved accuracy of the large time behavior of \(1\).](#) In particular, we consider open Newton-Cotes methods with  $p = 2, 4, 6$  points and we will also test a Gauss-Legendre quadrature. For the Gaussian quadrature we considered 8 points in each numerical cell. In the sequel, we will adopt the notation  $SP_k$ , with  $k = 2, 4, 6, G$ , to denote the SP schemes with (26) that is evaluated with second, fourth, sixth order Newton-Cotes quadrature or Gaussian quadrature, respectively. We highlight how possible singularities at the boundaries are avoided using open quadrature rules.

### 5.1 Test 1. Validation

In this subsection we consider a distribution function  $f(w, t)$ ,  $w \in [-1, 1] \times [-1, 1]$ , whose evolution is given by (1) in which, given the diffusion matrix  $\mathbb{D}$ , we chose the drift operator in such a way that the flux vanishes. In particular, we consider a linear drift term in the form (31) with a stationary state in the form

$$f_\infty(w) = C \exp\{-\phi(w)\}, \quad (40)$$

where  $\phi(w)$  is a given function of the state variable,  $C > 0$  a normalization constant. Therefore the linear drift term will be in the form

$$B(w) := -\nabla_w \cdot \mathbb{D}(w) - \mathbb{D}(w) \nabla_w \phi(w).$$

This is the case in which we have entropy dissipation and convergence of order  $p$ . In particular, we shall consider  $\mathbb{D}(w)$  a  $2 \times 2$  matrix of the form

$$\mathbb{D} = \begin{bmatrix} \frac{\sigma_1^2}{2}(1 - w_x^2)^2 & \rho \frac{\sigma_1 \sigma_2}{4}(1 - w_x^2)(1 - w_y^2) \\ \rho \frac{\sigma_1 \sigma_2}{4}(1 - w_x^2)(1 - w_y^2) & \frac{\sigma_2^2}{2}(1 - w_y^2)^2 \end{bmatrix}, \quad w_x, w_y \in [-1, 1]. \quad (41)$$

As initial condition we consider

$$f_0(w) = \beta [\exp(-c(w_x + \mu)^2) \exp(-c(w_y + \mu)^2) + \exp(-c(w_x - \mu)^2) \exp(-c(w_y - \mu)^2)], \quad (42)$$

with  $\mu = \frac{1}{2}$ ,  $c = 30$  and where  $\beta > 0$  is a normalization constant.

In Figure 1 we compute the relative  $L^1$  error of the numerical solution with respect to the exact stationary state  $f^\infty$  given by (40), i.e.

$$e_r^N(t^n) = \frac{\|f_{i,j}^n - f^\infty(w_{i,j})\|_{L^1}}{\|f^\infty(w_{i,j})\|_{L^1}} \quad (43)$$

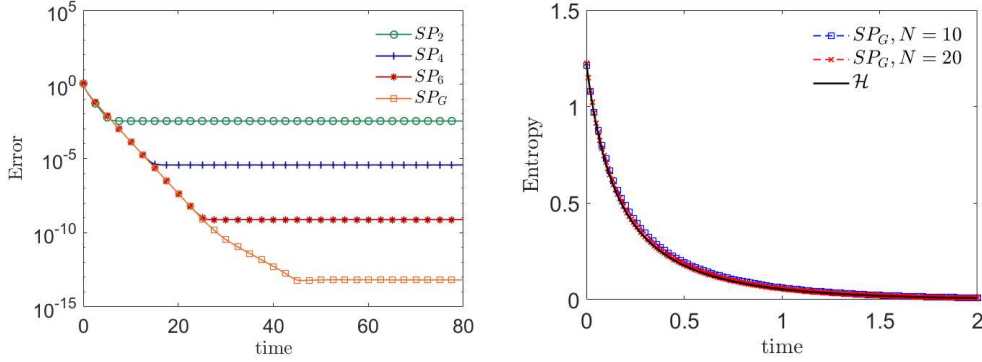


Figure 1: **Test 1.** Left: evolution over the time interval  $[0, 80]$  of the relative  $L^1$  error (43) computed with respect to the stationary solution (40) with  $\phi(w) = -d(w_x^8 + w_y^8)$ , where  $d = 12.5$ , for the  $SP_k$  scheme with different quadrature methods. Initial distribution as in (42) with  $\sigma_1^2 = \sigma_2^2 = 1$  and  $\rho = 0.9$ . We considered  $\Delta t = \Delta w / (20\sigma_1^2)$ ,  $\Delta w = 2/(N-1)$  and  $N = 81$ . Right: dissipation of the numerical entropy for  $SP_k$  scheme with Gaussian quadrature for two coarse grids with  $N = 10$  and  $N = 20$  points.

using  $N = 81$  grid points for the  $SP_k$  scheme with various quadrature rules. The different integration methods capture the steady state with different accuracy. In particular low order quadrature rules achieve their numerical steady state faster due to a saturation effect, whereas high order quadratures essentially reach machine precision in finite time. We considered in this plot semi-implicit time integration. In the same figure we illustrate how  $SP_k$  scheme dissipates the relative entropy (36) in the case of two coarse grids with  $N = 10$  and  $N = 20$  points.

In Table 1 we estimate the order of convergence of the schemes for first order time integration and a fourth order Runge-Kutta integration that is computed as  $\log_2(e_r^N(T))$ , with  $N = 81$  and  $T$  is the final time of the numerical test. The time step is chosen such that the CFL condition for the positivity of the scheme is satisfied, *i.e.*,  $\Delta t = \mathcal{O}(\Delta w^2)$ . We may observe that in the transient regime the second order is maintained, whilst we reach higher orders for large times, expressing the order of the quadrature rules. In Table 2 we estimate the order of convergence with first and second order semi-implicit methods. We chose the time step  $\Delta t = \mathcal{O}(\Delta w)$  to meet the positivity bound derived in Proposition 7. Here again, we may observe that the scheme is second order accurate in the transient regime and describes the long time behaviour of the problem with the order employed for the evaluation of the nonlinear weights.

## 5.2 Test 2. Alignment dynamics in bounded domains

In this test we provide numerical evidence of the failure of the scheme on models with non-vanishing flux at equilibrium. Let us consider the evolution of a distribution function as in (1) with  $w \in [-1, 1] \times [-1, 1]$ , anisotropic diffusion introduced in (41), and

$$\mathcal{B}[f](w, t) = \int_{[-1, 1] \times [-1, 1]} P(w, w_*) (w - w_*) f(w_*, t) dw_*, \quad (44)$$

with initial distribution of the form (42). We note that in this case we have no guarantee that the flux vanishes for large times.

First of all we consider (44) with  $P \equiv 1$ . This corresponds to  $\mathcal{B}[f](w, t) = B(w) = (w - U)$ ,  $U = \int_{\Omega} f(w_*, t) w_* dw_*$  that is constant since the mean of  $f$  is conserved. It is worth to notice that we are in the case in which a stationary state making the flux vanish may exist and we have decay of entropy. Since the stationary state of the problem is not known analytically, we computed the relative  $L^1$  error for successive approximations. We denote with  $f^{N_s}$  the approximation of  $f$  done using a grid with  $N_s$  points and we compute the error by considering as reference solution the one

|              |        |        |        |        |        |        |        |        |
|--------------|--------|--------|--------|--------|--------|--------|--------|--------|
| $\rho = 0.1$ | $SP_k$ |        |        |        | $SP_k$ |        |        |        |
| Time         | 2      | 4      | 6      | G      | 2      | 4      | 6      | G      |
| 1            | 1.9625 | 1.4962 | 1.6460 | 1.6461 | 1.9629 | 1.7472 | 1.8889 | 1.8891 |
| 10           | 1.9662 | 3.9708 | 7.3407 | 7.9144 | 1.9662 | 3.9708 | 7.4765 | 7.8903 |
| 20           | 1.9662 | 3.9708 | 7.4769 | 7.9144 | 1.9662 | 3.9708 | 7.4772 | 8.1457 |
| $\rho = 0.9$ | $SP_k$ |        |        |        | $SP_k$ |        |        |        |
| Time         | 2      | 4      | 6      | G      | 2      | 4      | 6      | G      |
| 1            | 1.8570 | 1.9049 | 1.9100 | 1.9100 | 1.8878 | 1.9559 | 1.9622 | 1.9622 |
| 10           | 1.9621 | 3.9678 | 2.1457 | 2.1554 | 1.9621 | 4.0880 | 2.4631 | 7.4904 |
| 20           | 1.9621 | 3.9800 | 6.0613 | 7.2470 | 1.9621 | 3.9800 | 6.0649 | 7.2697 |
| 50           | 1.9621 | 3.9800 | 6.2146 | 7.8973 | 1.9621 | 3.9800 | 6.2144 | 7.8964 |

Table 2: **Test 1.** Estimation of the order of convergence for  $SP_k$  scheme with first (left) and second order (right) semi-implicit methods. Rates have been computed using  $N = 21, 41, 81$  grid points,  $\sigma_1^2 = \sigma_2^2 = 1$ ,  $\Delta t = \Delta w / (20\sigma_1^2)$ , and two correlation coefficients  $\rho = 0.1$  (top) and  $\rho = 0.9$  (bottom).

|      |        |        |        |        |        |        |        |        |
|------|--------|--------|--------|--------|--------|--------|--------|--------|
|      | $SP_k$ |        |        |        | $SP_k$ |        |        |        |
| Time | 2      | 4      | 6      | G      | 2      | 4      | 6      | G      |
| 1    | 2.0830 | 2.1102 | 2.3204 | 2.4229 | 2.1320 | 2.3606 | 2.3602 | 2.3602 |
| 10   | 2.0914 | 2.2000 | 2.3614 | 2.5143 | 2.4199 | 2.8006 | 2.8195 | 2.8199 |
| 20   | 2.0914 | 3.7579 | 4.0746 | 3.8000 | 2.8741 | 3.7503 | 3.9163 | 3.8875 |

Table 3: **Test 2.** Estimation of the order of convergence for  $SP_k$  scheme with explicit Euler (left) and RK4 (right). Rates have been computed using  $N = 21, 41, 81$  grid points in each component of the computational cell. We considered  $P \equiv 1$ ,  $\sigma_1^2 = \sigma_2^2 = 1$ ,  $\rho = 0.1$ ,  $\Delta t = \Delta w^2 / (10\sigma_1^2 \Delta w + 10)$ .

of the successive refinement of the computational grid, i.e.

$$e_s(t^n) = \frac{\|f_{i,j}^{N_s,n} - f_{i,j}^{N_{s-1},n}\|_{L^1}}{\|f_{i,j}^{N_{s+1},n} - f_{i,j}^{N_s,n}\|_{L^1}}. \quad (45)$$

In detail, we chose  $N_1 = 21$ ,  $N_2 = 41$ ,  $N_3 = 81$  grid points. The order of convergence is then computed as  $\log_2(e_s(t^n))$ . In Table 3 we estimate the order of convergence of the  $SP_k$  scheme with explicit time integration methods. In particular, we present the case of first order forward Euler method and fourth order Runge-Kutta with a suitable time step to guarantee positivity of the scheme, i.e.  $\Delta t = O(\Delta w^2)$ . In Table 4 we estimate the order of convergence of the method in the case of semi-implicit time integration taking into account first and second order semi-implicit methods with  $\Delta t = O(\Delta w)$ . We may observe that in both cases, the proposed scheme is not capable to approximate the large time solution of the problem with high order. Indeed we have no theoretical guarantee that the flux vanishes at the equilibrium that is the assumption under which the scheme has been derived. The time evolution of the approximated solution are represented in Figure 2.

## Conclusion

We studied the construction of structure preserving methods for a class of **two-dimensional** Fokker-Planck equations with full nonconstant diffusion matrix. We have derived the schemes for stationary states that make the flux vanish. We have proved that mass conservation and positivity of the solution both with explicit and semi-implicit time integration hold even for problems with a non-vanishing flux at the steady state. Furthermore, the methods here developed are positivity

|      | $SP_k$ |        |        |        | $SP_k$ |        |        |        |
|------|--------|--------|--------|--------|--------|--------|--------|--------|
| Time | 2      | 4      | 6      | G      | 2      | 4      | 6      | G      |
| 1    | 1.9585 | 2.0242 | 2.2398 | 2.2615 | 1.9612 | 2.1190 | 2.2398 | 2.2732 |
| 10   | 2.0694 | 3.9977 | 3.6949 | 3.6477 | 2.0685 | 3.9643 | 3.6601 | 3.6140 |
| 20   | 2.0695 | 3.9982 | 3.6957 | 3.6486 | 2.0686 | 3.9643 | 3.6608 | 3.6140 |

Table 4: **Test 2.** Estimation of the order of convergence for  $SP_k$  scheme with first (left) and second order (right) semi-implicit integration. Rates have been computed using  $N = 21, 41, 81$ ,  $P \equiv 1$ ,  $\sigma_1^2 = \sigma_2^2 = 1$ ,  $\rho = 0.1$ ,  $\Delta t = \Delta w/(20\sigma_1^2)$ .

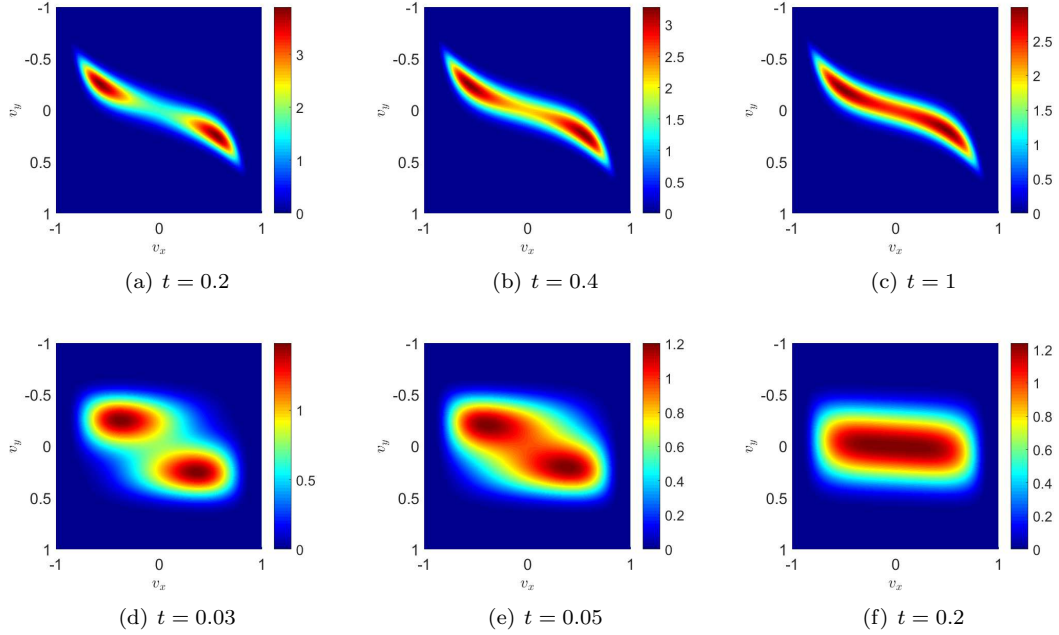


Figure 2: **Test 2.** Evolution of the numerical solution of the nonlinear Fokker-Planck equation with drift (44),  $P \equiv 1$ , and anisotropic diffusion matrix (41) with  $\sigma_1^2 = 0.1$ ,  $\sigma_2^2 = 0.5$  and correlation coefficient  $\rho = 0.9$  (top row) and  $\rho = 0.1$  (bottom row). The numerical solution has been computed with  $N = 101$  grid points in both components and semi-implicit time integration with  $\Delta t = \Delta w/(20 \max\{\sigma_1^2, \sigma_2^2\})$ .

preserving without any restriction on the discretization of the state variable both in the case of explicit and of semi-implicit time integration methods, the latter in particular lead to more mild restrictions on the time step that are very useful in the high-dimensional case. [On the other hand, the evolution scheme is in general equilibrium preserving for Fokker-Planck-type equations with nonconstant diffusion matrix if the drift is such that the flux function vanishes at the steady state and if it does not depend on the solution.](#) Under these assumptions we also showed entropy decay of the problem and that the introduced scheme dissipates the numerical entropy. We also presented numerical evidence of the theoretical findings. [Extensions of the proposed scheme to dimensions higher than two are currently under study and will be presented elsewhere.](#)

## Acknowledgements

This research was partially supported by the Italian Ministry of Education, University and Research (MIUR) through the “Dipartimenti di Eccellenza” Programme (2018-2022) – Department of Mathematical Sciences “G. L. Lagrange”, Politecnico di Torino (CUP: E11G18000350001).

Both authors are members of GNFM (Gruppo Nazionale per la Fisica Matematica) of INdAM (Istituto Nazionale di Alta Matematica), Italy.

NL would like to thank Compagnia San Paolo for financing her PhD scholarship.

## References

- [1] G. Albi, L. Pareschi, G. Toscani, and M. Zanella. Recent advances in opinion modeling: control and social influence. In N. Bellomo, P. Degond, and E. Tadmor, editors, *Active Particles Volume 1, Theory, Methods, and Applications*, Modeling and Simulation in Science, Engineering and Technology. Birkhäuser, 2017.
- [2] R. Bailo, J. A. Carrillo, and J. Hu. Fully discrete positivity-preserving and energy-decaying schemes for aggregation-diffusion equations with a gradient flow structure. Preprint [arxiv:1811.11502](#), 2018.
- [3] J. Barré, P. Degond, and E. Zatorska. Kinetic theory of particle interactions mediated by dynamical networks, *Multiscale Model. Simul.*, **15**(3): 1294–1323, 2017.
- [4] M. Bessemoulin-Chatard, and F. Filbet. A finite volume scheme for nonlinear degenerate parabolic equations, *SIAM J. Sci. Comput.*, **34**: 559–582, 2012.
- [5] M. Bessemoulin-Chatard, M. Herda, and T. Rey. Hypocoercivity and diffusion limit of a finite volume scheme for linear kinetic equations. Preprint [arXiv:1812.05967](#), 2018.
- [6] F. Bolley, J. A. Cañizo, and J. A. Carrillo. Stochastic mean-field limit: non-Lipschitz forces and swarming, *Math. Mod. Meth. Appl. Sci.*, **21**: 2179–2210, 2011.
- [7] S. Boscarino, F. Filbet, and G. Russo. High order semi-implicit schemes for time dependent partial differential equations, *J. Sci. Comput.*, **68**: 975–1001, 2016.
- [8] C. Buet, S. Cordier, and V. Dos Santos. A conservative and entropy scheme for a simplified model of granular media, *Transp. Theory Stat. Phys.*, **33**(2): 125–155, 2004.
- [9] C. Buet, and S. Dellacherie. On the Chang and Cooper numerical scheme applied to a linear Fokker-Planck equations, *Commun. Math. Sci.*, **8**(4): 1079–1090, 2010.
- [10] J. A. Carrillo, A. Chertock, Y. Huang. A finite-volume method for nonlinear non local equations with a gradient flow structure, *Commun. Comput. Phys.*, **17**(1): 233–258, 2015.
- [11] J. A. Carrillo, Y.-P. Choi, and L. Pareschi. Structure preserving schemes for the continuum Kuramoto model: phase transitions, *J. Comput. Phys.*, **376**: 365–389, 2019.

- [12] J. A. Carrillo, M. Fornasier, G. Toscani, and F. Vecil. Particle, kinetic, and hydrodynamic models of swarming. In G. Naldi, L. Pareschi, G. Toscani, editors, *Mathematical Modeling of Collective Behavior in Socio-Economic and Life Sciences*, Modeling and Simulation in Science, Engineering and Technology, Birkhäuser Boston, pp. 297–336, 2010.
- [13] C. Chainais-Hillairet, A. Jüngel, S. Schuchnigg. Entropy-dissipative discretization of nonlinear diffusion equations and discrete Beckner inequalities, *ESAIM Math. Model. Numer. Anal.*, **50**(1): 135–162, 2016.
- [14] J. S. Chang, and G. Cooper. A practical difference scheme for Fokker-Planck equations, *J. Comput. Phys.*, **6**(1): 1–16, 1970.
- [15] G. Dahlquist, and A. Björck. *Numerical Methods in Scientific Computing, Volume I*. SIAM, 2008.
- [16] A. Chauviere, T. Hillen and L. Preziosi. Modeling cell movement in anisotropic and heterogeneous network tissues, *Networks and Heterogeneous Media* **2** (2), 333–357, 2007.
- [17] G. Dimarco, and L. Pareschi. Numerical methods for kinetic equations, *Acta Numerica*, **23**: 369–520, 2014.
- [18] G. Dimarco, L. Pareschi, and M. Zanella. Uncertainty quantification for kinetic models in socio-economic and life sciences. In S. Jin, L. Pareschi, editors, *Uncertainty Quantification for Hyperbolic and Kinetic Equations*, SEMA SIMAI Springer Series, vol. 14, pp. 151–191, 2017.
- [19] R. Duan, M. Fornasier, and G. Toscani. A kinetic flocking model with diffusion, *Commun. Math. Phys.*, **300**: 95–145, 2010.
- [20] F. Filbet, L. Pareschi, T. Rey. On steady-state preserving spectral methods for homogeneous Boltzmann equations, *C R Acad Sci Paris, Ser-I*, **353** (4): 309–314, 2015
- [21] G. Furioli, A. Pulvirenti, E. Terraneo, and G. Toscani. Fokker–Planck equations in the modeling of socio-economic phenomena, *Math. Mod. Meth. Appl. Sci.*, **27**(1): 115–158, 2017.
- [22] L. Gosse. *Computing Qualitatively Correct Approximations of Balance Laws.*, SEMA SIMAI Springer Series, Springer, Berlin, 2013.
- [23] S. Gottlieb, C.-W. Shu, and E. Tadmor. Strong stability-preserving high-order time discretization methods, *SIAM Rev.*, **43**(1): 89–112, 2001.
- [24] R. Hegselmann, and U. Krause. Opinion dynamics and bounded confidence: Models, analysis, and simulation, *J. Artif. Soc. Soc. Simulat.*, **5**(3):1–33, 2002.
- [25] E. W. Larsen, C. D. Levermore, G. C. Pomraning, and J. G. Sanderson. Discretization methods for one-dimensional Fokker-Planck operators, *J. Comput. Phys.*, **61**(3): 359–390, 1985.
- [26] N. Loy, and L. Preziosi. Kinetic models with non-local sensing determining cell polarization and speed according to independent cues. *J. Math. Bio.* **80**, 373–421, 2020.
- [27] N. Loy, and A. Tosin. Markov jump processes and collision-like models in the kinetic description of multi-agent systems. *Comm. Math. Sci.* In press, 2020.
- [28] D. Matthes, A. Jüngel, and G. Toscani. Convex Sobolev inequalities derived from entropy dissipation, *Arch. Rat. Mech. Anal.*, **199**(2): 563–596, 2011.
- [29] M. Mohammadi, and A. Borzì. Analysis of the Chang–Cooper discretization scheme for a class of Fokker-Planck equations, *J. Numer. Math.* **23**(3):271–288, 2015



- [30] A. Okubo, and S.A. Levin. Diffusion and Ecological Problems: Modern Perspectives, *Springer, New York*, 2002 .
- [31] L. Pareschi, and T. Rey. Residual equilibrium schemes for time dependent partial differential equations, *Comput. Fluids*, **156**: 329–342, 2017.
- [32] L. Pareschi, G. Toscani. *Interacting Multiagent Systems: Kinetic equations and Monte Carlo methods*, Oxford University Press, 2013.
- [33] L. Pareschi, G. Toscani, A. Tosin, and M. Zanella. Hydrodynamic models of preference formation in multi-agent societies. *J. Nonlinear Sci.*, to appear.
- [34] L. Pareschi and M. Zanella. Structure preserving schemes for nonlinear Fokker-Planck equations and applications, *J. Sci. Comput.*, **74**(3):1575-1600, 2018.
- [35] L. Pareschi, and M. Zanella. Structure preserving schemes for mean-field equations of collective behavior. In M. Westdickenberg, C. Klingenberg, editors, *Theory, Numerics and Applications of Hyperbolic Problems II. HYP2016*, vol. 237 of *Springer Proceedings in Mathematics & Statistics*, pp. 405–421, Springer, Cham, 2018
- [36] Y. Qian, Z. Wang, and S. Zhou. A conservative, free energy dissipating, and positivity preserving finite difference scheme for multi-dimensional non local Fokker-Planck equation, *J. Comput. Phys.*, **386**: 22–36, 2019.
- [37] H. Risken. *The Fokker-Planck Equation, Methods of solution and Applications*, Springer-Verlag. Berlin, 1996.
- [38] D. L. Scharfetter, and H. K. Gummel. Large-signal analysis of a silicon Read diode oscillator, *IEEE Trans. Electron Devices*, **16**(1): 64–77, 1969.
- [39] P. C. da Silva, L. R. da Silva, E. K. Lenzi, R. S. Mendes, and L. C. Malacarne. Anomalous diffusion and anisotropic nonlinear Fokker-Planck equation, *Physica A*, **342**: 16–21, 2004.
- [40] G. Toscani. Entropy production and the rate of convergence to equilibrium for the Fokker-Planck equation, *Quart. Appl. Math.*, **57**: 521–541, 1999.
- [41] G. Toscani. Kinetic models of opinion formation, *Commun. Math. Sci.*, **4**(3): 481–496, 2006.
- [42] G. Toscani, and C. Villani. Sharp entropy dissipation bounds and explicit rate of trend to equilibrium for the spatially homogeneous Boltzmann equation, *Commun. Math. Phys.*, **203**(3): 667–706, 1999.
- [43] A. Tosin, and M. Zanella. Boltzmann-type models with uncertain binary interactions, *Commun. Math. Sci.*, **16**(4): 962-984, 2018.
- [44] C. Yates, R. Erban, C. Escudero, L. Couzin, J. Buhl, L. Kevrekidis, P. Maini and D. Sumpter. Inherent noise can facilitate coherence in collective swarm motion *Proc. Nat. Acad. Sci.*, **106**(14): 5464–5469, 2009.