



**HAL**  
open science

# On the square-root approximation finite volume scheme for nonlinear drift-diffusion equations

Clément Cancès, Juliette Venel

► **To cite this version:**

Clément Cancès, Juliette Venel. On the square-root approximation finite volume scheme for nonlinear drift-diffusion equations. 2022. hal-03693887v1

**HAL Id: hal-03693887**

**<https://hal.science/hal-03693887v1>**

Preprint submitted on 13 Jun 2022 (v1), last revised 12 Sep 2022 (v3)

**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.

# ON THE SQUARE-ROOT APPROXIMATION FINITE VOLUME SCHEME FOR NONLINEAR DRIFT-DIFFUSION EQUATIONS

CLÉMENT CANCÈS AND JULIETTE VENEL

ABSTRACT. We study a finite volume scheme for the approximation of the solution to convection diffusion equations with nonlinear convection and Robin boundary conditions. The scheme builds on the interpretation of such a continuous equation as the hydrodynamic limit of some simple exclusion jump process. We show that the scheme admits a unique discrete solution, that the natural bounds on the solution are preserved, and that it encodes the second principle of thermodynamics in the sense that some free energy is dissipated along time. The convergence of the scheme is then rigorously established thanks to compactness arguments. Numerical simulations are finally provided, highlighting the overall good behavior of the scheme.

## 1. PRESENTATION OF THE PROBLEM

1.1. **The governing equations.** In this paper, we focus on the simple yet already interesting nonlinear Fokker-Planck equation

$$(1a) \quad \partial_t \rho + \nabla \cdot F = 0,$$

$$(1b) \quad F + \eta(\rho) \nabla \phi + \nabla \rho = 0,$$

set on a connected bounded open subset  $\Omega$  of  $\mathbb{R}^d$ , which is further assumed to be polyhedral in what follows, and for positive times  $t \geq 0$ . Its (finite) Lebesgue measure is denoted by  $m_\Omega$ . While diffusion is linear, convection is not since one considers a degenerate mobility function  $\eta$  of the form

$$(2) \quad \eta(\rho) = \rho(1 - \rho)$$

accounting for volume-filling to enforce  $0 \leq \rho \leq 1$ . In (1), the potential  $\phi \in W^{1,\infty}(\Omega)$  (referred as the electric potential in what follows) is assumed to be given, and nonnegative without loss of generality:  $\phi \geq 0$ . Our purpose can be extended to the case of a self-consistent electric potential  $\phi$  related to the charge density  $\rho$  through a Poisson equation without other difficulties than those that are already addressed in the literature, see for instance [8].

The system we consider is not isolated as in [6], but rather in interaction with a surrounding environment through its boundary  $\Gamma = \partial\Omega$ . More precisely, we assume that there exist  $\alpha, \beta \in W^{1,\infty}(\Gamma)$  with  $\alpha(x) > \beta(x) > 0$  for all  $x \in \Gamma$  such that

$$(3) \quad F \cdot \nu = \alpha\rho - \beta \quad \text{on } \mathbb{R}_+ \times \Gamma,$$

where  $\nu$  denotes the normal to  $\Gamma$  outward w.r.t.  $\Omega$ . The system is complemented by an initial condition  $\rho^0$  compatible with the volume-filling constraint:

$$(4) \quad \rho|_{t=0} = \rho^0 \in L^\infty(\Omega; [0, 1]).$$

---

2000 *Mathematics Subject Classification.* 65M08, 65M12, 35K51, 35Q92, 92D25.

*Key words and phrases.* Nonlinear convection diffusion, finite volume method, energy dissipation, convergence.

Our goal is to provide some provably convergent approximation of the problem (1)–(4). The stability of our numerical method, to be detailed in Section 2, mimics some stability features of the continuous problem inherited from thermodynamics.

**1.2. Energy dissipation structure.** The system (1)–(4) under consideration inherits some key property from thermodynamics. Defining its free energy by

$$\mathcal{F}(\rho) = \int_{\Omega} (h(\rho) + \rho\phi), \quad h(\rho) = \rho \log(\rho) + (1 - \rho) \log(1 - \rho) + \log(2) \geq 0,$$

then it is dissipated within  $\Omega$ , but energy coming from the surrounding environment can enter the system thanks to the boundary flux (3).

Introducing the chemical and electrochemical potentials  $\mu$  and  $\xi$  respectively defined by

$$(5) \quad \mu = h'(\rho) = \log \frac{\rho}{1 - \rho}, \quad \xi = \mu + \phi = \frac{\delta \mathcal{F}}{\delta \rho}(\rho), \quad 0 < \rho < 1,$$

the chain rule  $\nabla \rho = \eta(\rho) \nabla \mu$  allows to reformulate the flux

$$(6) \quad F = -\eta(\rho) \nabla(\phi + \mu) = -\eta(\rho) \nabla \xi.$$

On the other hand, setting

$$(7) \quad \xi^{\Gamma} = \phi - \log(\alpha/\beta - 1) \in W^{1,\infty}(\Gamma) \quad \text{and} \quad \kappa = \sqrt{\beta(\alpha - \beta)} \in W^{1,\infty}(\Gamma),$$

the boundary flux (3) can be expressed by the mean of a Butler-Volmer type formula:

$$(8) \quad F \cdot \nu = \kappa \left( \rho e^{\frac{1}{2}(\phi - \xi^{\Gamma})} - (1 - \rho) e^{-\frac{1}{2}(\phi - \xi^{\Gamma})} \right) = 2\kappa \sqrt{\rho(1 - \rho)} \sinh \left( \frac{1}{2}(\xi - \xi^{\Gamma}) \right).$$

The quantity  $\xi^{\Gamma}$  has to be thought as an electrochemical potential associated to the surrounding environment. When a quantity  $n^{\Gamma} = \int_{\Gamma} F \cdot \nu$  of the chemical species of interest enters (resp. leaves)  $\Omega$ , the income (resp. loss) in free energy is equal to  $n^{\Gamma} \xi^{\Gamma}$ . Therefore, the total free energy defined (up to an additive constant) by

$$(9) \quad \mathcal{F}_{\text{tot}}(t) = \mathcal{F}(\rho(t)) + \int_0^t \int_{\Gamma} \xi^{\Gamma} F \cdot \nu, \quad t \geq 0,$$

corresponds to the whole isolated system made of  $\Omega$  and its surrounding environment. As the following proposition shows, it is decaying along time.

**Proposition 1.1.** *Let  $\rho$  be a strong solution to (1)–(4), then*

$$(10) \quad \mathcal{F}_{\text{tot}}(t) \leq \mathcal{F}_{\text{tot}}(s) \leq \mathcal{F}(\rho^0) \leq (\|\phi\|_{\infty} + \log(2)) m_{\Omega}, \quad t \geq s \geq 0.$$

*Moreover, there exists  $C_1$  depending on  $\Gamma, \alpha, \beta$  and  $\phi$  such that*

$$(11) \quad \mathcal{F}_{\text{tot}}(t) \geq -C_1 t, \quad t \geq 0.$$

*Proof.* On first remarks that thanks to its definition (9), the initial total free energy coincides with the free energy contained in  $\Omega$ , i.e.  $\mathcal{F}_{\text{tot}}(0) = \mathcal{F}(\rho^0)$  thanks to (4). The bound on the initial energy  $\mathcal{F}(\rho^0)$  is readily deduced from  $0 \leq \rho^0 \leq 1$  and  $0 \leq h(\rho^0) \leq \log(2)$ . Let us now check that  $\mathcal{F}_{\text{tot}}$  is decaying along time. To this end, let us compute

$$(12) \quad \frac{d\mathcal{F}_{\text{tot}}}{dt}(t) = \int_{\Omega} \xi \partial_t \rho + \int_{\Gamma} \xi^{\Gamma} F \cdot \nu = \int_{\Omega} F \cdot \nabla \xi + \int_{\Gamma} (\xi^{\Gamma} - \xi) F \cdot \nu.$$

Both terms on the right-hand side are nonpositive respectively because of (6) and (8), so that (10) holds true.

To establish (11), one only has to notice that  $\mathcal{F}(\rho(t))$  is nonnegative for all  $t \geq 0$ , so that

$$\mathcal{F}_{\text{tot}}(t) = \mathcal{F}(\rho(t)) + \int_0^t \int \xi^\Gamma F \cdot \nu \geq \int_0^t \int \xi^\Gamma F \cdot \nu \geq -t \|\xi^\Gamma\|_\infty \|F \cdot \nu\|_\infty, \quad t \geq 0.$$

Uniform bounds on  $\xi^\Gamma$  and on  $F \cdot \nu$  easily follow from their expressions (7) and (3) together with  $0 \leq \rho \leq 1$ .  $\square$

The estimates highlighted in Proposition 1.1 encode some strong stability in the system (1)–(4). The precise quantification of the dissipation rate of the total free energy even provides sufficiently compactness to establish the existence of weak solutions to (1)–(4). The numerical method we introduce in the next section satisfies similar energy dissipation estimates, on which the numerical analysis we propose relies.

**Definition 1.2.** *A function  $\rho$  is said to be a weak solution to (1)–(4) if:*

- (i)  $\rho$  belongs to  $L^\infty(\mathbb{R}_+ \times \Omega; [0, 1]) \cap L^2_{loc}(\mathbb{R}_+; H^1(\Omega))$ , hence its trace  $\gamma\rho$  on  $\mathbb{R}_+ \times \Gamma$  belongs to  $L^\infty(\mathbb{R}_+ \times \Gamma; [0, 1]) \cap L^2_{loc}(\mathbb{R}_+; H^{1/2}(\Gamma))$ ;
- (ii) for all  $\varphi \in C_c^\infty([0, T] \times \overline{\Omega})$ , the following equality holds:

$$(13) \quad \iint_{\mathbb{R}_+ \times \Omega} \rho \partial_t \varphi + \int_\Omega \rho^0 \varphi(0, \cdot) - \iint_{\mathbb{R}_+ \times \Omega} (\eta(\rho) \nabla \phi + \nabla \rho) \cdot \nabla \varphi - \iint_{\mathbb{R}_+ \times \Gamma} (\alpha \gamma \rho + \beta) \varphi = 0.$$

**1.3. Goal and positioning of the paper.** The goal of this paper is to propose a seemingly new scheme to approximate nonlinear drift diffusion equations of the form (1). Such nonlinear drift diffusion problem arises in many contexts that are often more complex than the simple one prescribed by (1). We could for instance think about systems involving several species, coupled either via cross-diffusion [6], or via a self-consistent electric potential [7]. We claim that a large part of our work (in practice all excepted what is related to uniqueness) can be transposed to the more complex setting of [7]. To enlighten the presentation, we rather adopt here a simpler setting, where the potential  $\phi$  is given, but still with boundary conditions of Butler-Volmer type.

Even though this scheme has a very natural probabilistic interpretation in terms of jump process, its use with a deterministic approach to compute solutions to (1) has not been explored so far up to our knowledge. The scheme can be thought as an extension to the case of a nonlinear mobility function  $\eta$  defined by (2) of the approach proposed by [22] and studied in [17], even though the method proposed therein is mesh-less and yields non-explicit diffusion tensors at the limit we avoid here.

Our study covers several aspects. First, since our scheme is implicit, it yields a nonlinear system for which we show well-posedness and the preservation of the  $L^\infty$  bounds. These properties follow from the monotonicity of the scheme. Another interesting aspect of the scheme is its free energy stability: a discrete counterpart to Proposition 1.1 is established. Schemes encoding the second principle of thermodynamics have raised an important interest in the last years. In the case of a linear mobility  $\eta(\rho) = \rho$ , the Scharfetter-Gummel scheme [11], the SQRA scheme [17], or the Chang-Cooper scheme [5] are popular solutions since the scheme for solving the resulting linear Fokker-Planck equation amounts to the resolution of a linear system, in opposition to more involved strategies building on the Wasserstein gradient flow interpretation of the continuous problem (with no-flux boundary conditions), see for instance [23, 3, 9, 21, 10]. We are not aware of extensions to the nonlinear mobility framework (2) for the Scharfetter-Gummel or the Chang-Cooper scheme. In opposition, the extension of the SQRA scheme proposed in this paper is very natural.

Second, we mathematically assess the convergence of the scheme when the discretization parameters (mesh size and time step) tend to 0. To this end, one needs to properly quantify the free energy dissipation. This is done thanks to primal and dual dissipation potentials inspired from [24, 25]. The convergence proof then relies on compactness arguments, following the strategy of [15]. Our convergence result is not quantitative, since no error estimate has been derived so far. Then we show in the numerical experiments that the scheme is second order accurate in space and first order in time. See for instance [18], where error estimates for several schemes including SQRA finite volumes are derived for steady linear Fokker-Planck equations. We also highlight the fact that the resolution of the nonlinear system by the Newton-Raphson method is very efficient, even for large CFL conditions. The only drawback we have noticed so far for our scheme is its loss of accuracy in the large Péclet regime.

## 2. THE FINITE VOLUME SCHEME AND MAIN RESULTS

Before introducing the so-called *square-root approximation* (SQRA) scheme, one first needs to introduce some notation related to space and time discretizations.

**2.1. Space and time discretizations of  $\mathbb{R}_+ \times \Omega$ .** The SQRA finite volume scheme enters the framework of two-point flux approximation (TPFA) finite volumes, which are known to yield very efficient schemes but require meshes fulfilling the well-known orthogonality condition (iii) below, see for instance [14, 16].

**Definition 2.1.** *An admissible mesh of  $\Omega$  is a triplet  $(\mathcal{T}, \mathcal{E}, (x_K)_{K \in \mathcal{T}})$  such that the following conditions are fulfilled.*

- (i) *Each control volume (or cell)  $K \in \mathcal{T}$  is non-empty, open, polyhedral and convex. We assume that*

$$K \cap L = \emptyset \quad \text{if } K, L \in \mathcal{T} \text{ with } K \neq L, \quad \text{while} \quad \bigcup_{K \in \mathcal{T}} \overline{K} = \overline{\Omega}.$$

- (ii) *Each face  $\sigma \in \mathcal{E}$  is closed and is contained in a hyperplane of  $\mathbb{R}^d$ , with positive  $(d-1)$ -dimensional Hausdorff (or Lebesgue) measure denoted by  $m_\sigma = \mathcal{H}^{d-1}(\sigma) > 0$ . We assume that  $\mathcal{H}^{d-1}(\sigma \cap \sigma') = 0$  for  $\sigma, \sigma' \in \mathcal{E}$  unless  $\sigma' = \sigma$ . For all  $K \in \mathcal{T}$ , we assume that there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that  $\partial K = \bigcup_{\sigma \in \mathcal{E}_K} \sigma$ . Moreover, we suppose that  $\bigcup_{K \in \mathcal{T}} \mathcal{E}_K = \mathcal{E}$ . Given two distinct control volumes  $K, L \in \mathcal{T}$ , the intersection  $\overline{K} \cap \overline{L}$  either reduces to a single face  $\sigma \in \mathcal{E}$  denoted by  $K|L$ , or its  $(d-1)$ -dimensional Hausdorff measure is 0.*
- (iii) *The cell-centers  $(x_K)_{K \in \mathcal{T}}$  are two by two distinct points of  $\Omega$ . If  $K, L \in \mathcal{T}$  share a face  $K|L$ , then the vector  $x_L - x_K$  is orthogonal to  $K|L$  and oriented from  $K$  to  $L$ .*
- (iv) *For the boundary faces  $\sigma \subset \partial\Omega$ , we assume that there exists  $x_\sigma \in \sigma$  such that  $x_\sigma - x_K$  is orthogonal to  $\sigma$ .*

In the above definition, we do not suppose that  $x_K$  belongs to  $K$ . We allow for more general grids, like for instance Delaunay triangulation or Laguerre cells. The condition on the fact that the  $x_K$  are two-by-two distinct is not restrictive: if two cell centers  $x_K$  and  $x_L$  coincide, one just has to merge the two cells  $K$  and  $L$  and to remove  $K|L$  from  $\mathcal{E}$ .

We denote by  $m_K$  the  $d$ -dimensional Lebesgue measure of the control volume  $K$ . The set of the faces is partitioned into two subsets: the set  $\mathcal{E}_{\text{int}}$  of the interior faces defined by

$$\mathcal{E}_{\text{int}} = \{\sigma \in \mathcal{E} \mid \sigma = K|L \text{ for some } K, L \in \mathcal{T}\},$$

and the set  $\mathcal{E}_{\text{ext}}$  of the exterior faces defined by

$$\mathcal{E}_{\text{ext}} = \{\sigma \in \mathcal{E} \mid \sigma \subset \partial\Omega\}.$$

For a given control volume  $K \in \mathcal{T}$ , we also define  $\mathcal{E}_{K,\text{int}} = \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$  and  $\mathcal{E}_{K,\text{ext}} = \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$  the sets of its internal and external faces. We may write  $\sigma = K|L$  to signify that  $\sigma \in \mathcal{E}_{K,\text{int}}$ . For such internal edges  $\sigma = K|L$ , we denote by  $x_\sigma$  the intersection between  $[x_K, x_L]$  and the hyperplane containing  $\sigma$ . Note that  $x_\sigma$  does not necessarily belong to  $\sigma$ .

In what follows, we denote by

$$d_\sigma = \begin{cases} |x_K - x_L| & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ |x_K - x_\sigma| & \text{if } \sigma \in \mathcal{E}_{\text{ext}}, \end{cases} \quad a_\sigma = \frac{m_\sigma}{d_\sigma}, \quad \sigma \in \mathcal{E}.$$

We also define the signed distance  $d_{K\sigma}$  between  $x_K$  and  $\sigma \in \mathcal{E}_K$  thanks to the relation

$$d_{K\sigma} \nu_{K\sigma} = x_\sigma - x_K, \quad \sigma \in \mathcal{E}_K, K \in \mathcal{T},$$

where  $\nu_{K\sigma}$  stands for the normal to  $\sigma$  outward w.r.t.  $K$ . Even though  $d_{K\sigma}$  can take negative values for interior faces, one still has

$$d_{K\sigma} + d_{L\sigma} = d_\sigma > 0 \quad \text{for } \sigma = K|L \in \mathcal{E}_{\text{int}},$$

as well as the geometric relation

$$m_K = \frac{1}{d} \sum_{\sigma \in \mathcal{E}_K} m_\sigma d_{K\sigma}, \quad K \in \mathcal{T}.$$

We further introduce the *size*  $\delta_{\mathcal{T}}$  and the *regularity factor*  $\zeta_{\mathcal{T}}$  of the mesh:

$$(14) \quad \delta_{\mathcal{T}} = \max_{K \in \mathcal{T}} \text{diam}(K), \quad \zeta_{\mathcal{T}} = \max_{K \in \mathcal{T}} \max_{\sigma \in \mathcal{E}_K} \left( \frac{\text{diam}(K)}{d_\sigma} + \frac{d_\sigma}{\text{diam}(K)} \right).$$

Given  $\mathbf{u} = ((u_K)_{K \in \mathcal{T}}, (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}) \in \mathbb{R}^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$ , then for all  $K \in \mathcal{T}$ , we define the mirror value of  $u_K$  w.r.t.  $\sigma \in \mathcal{E}_K$  by

$$(15) \quad u_{K\sigma} = \begin{cases} u_L & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ u_\sigma & \text{if } \sigma \in \mathcal{E}_{\text{ext}}. \end{cases}$$

Concerning the time discretization, we consider for notation simplicity a uniform time stepping. More precisely, a time discretization is given by the choice of a time step  $\tau > 0$ , from which we construct discrete times  $t_n = n\tau$ ,  $n \geq 0$ . We stress that our study can be extended without any particular difficulty to the case of non-uniform time discretizations.

**2.2. The SQRA finite volume scheme.** Given an admissible discretization  $(\mathcal{T}, \mathcal{E}, (x_K)_{K \in \mathcal{T}})$  of  $\Omega$  and a time step  $\tau$ , let us detail the scheme to be studied in this paper. First, the initial data  $\rho^0$  is discretized into  $\boldsymbol{\rho}^0 = (\rho_K^0)_{K \in \mathcal{T}} \in [0, 1]^{\mathcal{T}}$  by setting

$$(16) \quad \rho_K^0 = \frac{1}{m_K} \int_K \rho^0, \quad K \in \mathcal{T}.$$

The potential  $\phi$  is discretized into  $\boldsymbol{\phi} = ((\phi_K)_{K \in \mathcal{T}}; (\phi_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}})$  by setting

$$(17) \quad \phi_K = \phi(x_K) \quad \text{and} \quad \phi_\sigma = \phi(x_\sigma), \quad K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}.$$

As usual in the finite volume context, the conservation law (1a) is discretized into

$$(18) \quad \frac{\rho_K^n - \rho_K^{n-1}}{\tau} m_K + \sum_{\sigma \in \mathcal{E}_K} m_\sigma F_{K\sigma}^n = 0, \quad K \in \mathcal{T}, n \geq 1.$$

The index  $n$  for the numerical flux  $F_{K\sigma}^n$  across  $\sigma$  outward w.r.t.  $K$  in (18) indicates that our time discretization strategy relies on the backward Euler scheme. The bulk numerical fluxes are then defined by

$$(19a) \quad F_{K\sigma}^n = \frac{1}{d_\sigma} \left[ \rho_K^n (1 - \rho_L^n) e^{\frac{1}{2}(\phi_K - \phi_L)} - \rho_L^n (1 - \rho_K^n) e^{\frac{1}{2}(\phi_L - \phi_K)} \right] \quad \text{for } \sigma = K|L \in \mathcal{E}_{\text{int}}.$$

To preserve the second order accuracy in space, the boundary condition (3) is discretized by setting

$$(19b) \quad F_{K\sigma}^n = \frac{1}{d_\sigma} \left[ \rho_K^n (1 - \rho_\sigma^n) e^{\frac{1}{2}(\phi_K - \phi_\sigma)} - \rho_\sigma^n (1 - \rho_K^n) e^{\frac{1}{2}(\phi_\sigma - \phi_K)} \right] = \alpha_\sigma \rho_\sigma^n - \beta_\sigma, \quad \text{for } \sigma \in \mathcal{E}_{\text{ext}},$$

where, having set  $\alpha_\sigma = \alpha(x_\sigma)$  and  $\beta_\sigma = \beta(x_\sigma)$ ,

$$(20) \quad \rho_\sigma^n = \frac{d_\sigma \beta_\sigma + \rho_K^n e^{\frac{1}{2}(\phi_K - \phi_\sigma)}}{d_\sigma \alpha_\sigma + \rho_K^n e^{\frac{1}{2}(\phi_K - \phi_\sigma)} + (1 - \rho_K^n) e^{-\frac{1}{2}(\phi_K - \phi_\sigma)}}$$

is the unique value achieving the second equality in (19b). With a slight abuse of notation, we still denote by  $\boldsymbol{\rho}^n = ((\rho_K^n)_{K \in \mathcal{T}}, (\rho_\sigma^n)_{\sigma \in \mathcal{E}_{\text{ext}}})$  the discrete density enriched with its boundary edge values prescribed by (20).

Formula (19a) can be interpreted as a Butler-Volmer law located at the interface between the cells  $K$  and  $L$ . The probability that a particle jumps from  $K$  to  $L$  is proportional to the number  $\rho_K^n$  of candidates in  $K$  for a jump as well as to the number of available sites  $(1 - \rho_L^n)$  to host the particle in cell  $L$ . The drift  $\phi_K - \phi_L$  appears in an exponential with balanced prefactors  $1/2$ , which is natural since  $K$  and  $L$  play symmetric roles in the formula. The scheme (18)&(19) is then a simple backward Euler discretisation of the dynamics prescribed by the infinitesimal generator of a weakly asymmetric simple exclusion process (WASEP), see [19].

Assume now that  $\boldsymbol{\rho}^n \in (0, 1)^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$  (this will be rigorously established later on, see Lemma 3.1). The consistency of formula (19a) with (1b) follows from the identity

$$(21) \quad F_{K\sigma}^n = \frac{2}{d_\sigma} \eta_\sigma^n \sinh \left( \frac{\xi_K^n - \xi_{K\sigma}^n}{2} \right),$$

with  $\xi_K^n = h'(\rho_K^n) + \phi_K$  for  $K \in \mathcal{T}$ ,  $\xi_\sigma^n = h'(\rho_\sigma^n) + \phi_\sigma$  for  $\sigma \in \mathcal{E}_{\text{ext}}$ , and where  $\xi_{K\sigma}^n$  is the mirror value of  $\xi_K^n$  in the sense of (15). Moreover, we have set

$$(22) \quad \eta_\sigma^n = \sqrt{\rho_K^n (1 - \rho_K^n) \rho_{K\sigma}^n (1 - \rho_{K\sigma}^n)} = \sqrt{\eta(\rho_K^n) \eta(\rho_{K\sigma}^n)},$$

Taylor expanding formula (21), one gets that for each  $n \geq 1$  and  $\sigma \in \mathcal{E}$ , there exists  $r_\sigma^n \in (0, 1)$  such that

$$F_{K\sigma}^n = \frac{\eta_\sigma^n}{d_\sigma} \left( 2 \sinh \left( \frac{h'(\rho_K^n) - h'(\rho_{K\sigma}^n)}{2} \right) + (\phi_K - \phi_{K\sigma}) \cosh \left( \frac{h'(\rho_K^n) - h'(\rho_{K\sigma}^n)}{2} \right) + \frac{(\phi_K - \phi_{K\sigma})^2}{4} \sinh \left( \frac{h'(\rho_K^n) - h'(\rho_{K\sigma}^n) + r_\sigma^n (\phi_K - \phi_{K\sigma})}{2} \right) \right).$$

Then using the identities

$$(23) \quad \eta_\sigma^n \sinh\left(\frac{h'(\rho_K^n) - h'(\rho_{K\sigma}^n)}{2}\right) = \frac{\rho_K^n - \rho_{K\sigma}^n}{2},$$

$$(24) \quad \eta_\sigma^n \cosh\left(\frac{h'(\rho_K^n) - h'(\rho_{K\sigma}^n)}{2}\right) = \frac{\eta(\rho_K^n) + \eta(\rho_{K\sigma}^n)}{2} + \frac{(\rho_K^n - \rho_{K\sigma}^n)^2}{2},$$

and  $\sinh(a+b) = \sinh(a)\cosh(b) + \sinh(b)\cosh(a)$ , we get that

$$(25) \quad F_{K\sigma}^n = \frac{\rho_K^n - \rho_{K\sigma}^n}{d_\sigma} + \frac{\eta(\rho_K^n) + \eta(\rho_{K\sigma}^n)}{2} \frac{\phi_K - \phi_{K\sigma}}{d_\sigma} + R_\sigma^n,$$

with

$$(26) \quad R_\sigma^n = \frac{\phi_K - \phi_{K\sigma}}{2d_\sigma} (\rho_K^n - \rho_{K\sigma}^n)^2 + \frac{(\phi_K - \phi_{K\sigma})^2}{8} \frac{\rho_K^n - \rho_{K\sigma}^n}{d_\sigma} \cosh\left(\frac{r_\sigma^n(\phi_K - \phi_{K\sigma})}{2}\right) \\ + \frac{(\phi_K - \phi_{K\sigma})^2}{8d_\sigma} (\eta(\rho_K^n) + \eta(\rho_{K\sigma}^n) + (\rho_K^n - \rho_{K\sigma}^n)^2) \sinh\left(\frac{r_\sigma^n(\phi_K - \phi_{K\sigma})}{2}\right).$$

Let  $\bar{\rho} : \mathbb{R}_+ \times \bar{\Omega} \rightarrow (0, 1)$  be a smooth (say  $C^{0,1}$  in time and  $C^{1,1}$  in space) function, then for all  $n \geq 1$ , define  $\bar{\rho}_K^n = \bar{\rho}(t_n, x_K)$ ,  $K \in \mathcal{T}$ , and

$$\bar{F}_{K\sigma}^n = \frac{1}{d_\sigma} \left[ \bar{\rho}_K^n (1 - \bar{\rho}_L^n) e^{\frac{1}{2}(\phi_K - \phi_L)} - \bar{\rho}_L^n (1 - \bar{\rho}_K^n) e^{\frac{1}{2}(\phi_L - \phi_K)} \right], \quad \sigma = K|L \in \mathcal{E}_{\text{int}}.$$

In the case of a uniform cartesian grid, where  $x_\sigma = \frac{x_K + x_L}{2}$  is the center of mass of  $\sigma$ , it results from the expression (25) of the flux that

$$\bar{F}_{K\sigma}^n = \frac{1}{m_\sigma} \int_\sigma \bar{F}(t_n) \cdot \nu_{K\sigma} + \mathcal{O}(d_\sigma^2),$$

where  $\bar{F} = -\nabla \bar{\rho} - \eta(\bar{\rho}) \nabla \phi$  is the flux corresponding to  $\bar{\rho}$ . The SQRA scheme, which owes its name to the choice (22) of a geometric mean for the edge mobilities  $\eta_\sigma^n$  and to the fact that it extends to the nonlinear mobility setting the linear SQRA scheme [22, 17], is then expected to be second order accurate w.r.t. space and first order accurate w.r.t. time since it relies on the backward Euler scheme. This will be confirmed by the numerical results exhibited in Section 5.

**2.3. Our main results and organisation of the paper.** Even though finer results can be found in the Sections 3 and 4 devoted to their proofs, we state here simple presentations of our main results. The first one, namely Theorem 2.2, is related to the characteristics of the scheme given a fixed mesh  $(\mathcal{T}, \mathcal{E}, (x_K)_{K \in \mathcal{T}})$  and time step  $\tau$ . We show in particular that the scheme is well posed, preserves the  $L^\infty$  bounds and is free-energy diminishing, in the sense that the discrete solution satisfies a discrete counterpart of Proposition 1.1. Then Theorem 2.3 states the convergence of the approximate solution provided by the scheme (16)–(19) towards the weak solution to (1)–(4) as the size of the mesh  $\delta_\mathcal{T}$  and the time step  $\tau$  tend to 0. The convergence analysis strongly relies on the energy stability of the scheme, and more precisely on the quantification of the free energy dissipation.

Given  $\boldsymbol{\rho}^n = (\rho_K^n)_{K \in \mathcal{T}} \in [0, 1]^\mathcal{T}$ , then we define

$$(27) \quad \mathcal{F}_\mathcal{T}(\boldsymbol{\rho}^n) = \sum_{K \in \mathcal{T}} m_K (h(\rho_K^n) + \phi_K \rho_K^n), \quad \mathcal{F}_{\mathcal{T}, \text{tot}}^n = \mathcal{F}_\mathcal{T}(\boldsymbol{\rho}^n) + \sum_{p \geq 1} \tau \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m_\sigma \xi_\sigma^\Gamma F_{K\sigma}^p,$$



where the external fluxes  $F_{K\sigma}^p$  are related to  $\rho^p$  through formula (19b), and where, consistently with (7), we have set

$$(28) \quad \xi_\sigma^\Gamma = \phi_\sigma - \log\left(\frac{\alpha_\sigma}{\beta_\sigma} - 1\right), \quad \sigma \in \mathcal{E}_{\text{ext}}.$$

Initially, both energies coincide:  $\mathcal{F}_\mathcal{T}(\rho^0) = \mathcal{F}_{\mathcal{T},\text{tot}}^0$ , and it follows from Jensen's inequality and from the regularity of  $\phi$  that

$$(29) \quad \mathcal{F}_\mathcal{T}(\rho^0) \leq \mathcal{F}(\rho^0) + 2\|\nabla\phi\|_\infty \delta_\mathcal{T} m_\Omega.$$

In particular,  $\mathcal{F}_\mathcal{T}(\rho^0)$  is bounded uniformly w.r.t.  $\delta_\mathcal{T}$  owing to (10) and to  $\delta_\mathcal{T} \leq \text{diam}(\Omega)$ .

**Theorem 2.2.** *Given  $\rho^{n-1} \in [0, 1]^\mathcal{T}$ , there exists a unique solution  $\rho^n \in (0, 1)^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$  to the non-linear system corresponding to the scheme (16)–(20). Moreover,*

$$(30) \quad \mathcal{F}_{\mathcal{T},\text{tot}}^{n-1} \geq \mathcal{F}_{\mathcal{T},\text{tot}}^n \geq -C_1 t_n, \quad n \geq 1,$$

with  $C_1$  as in Proposition 1.1.

Theorem 2.2 is a partial presentation of the results established in Section 3. Interested readers can find there some precise quantification of the dissipated total free energy we do not mention here to keep the presentation simple.

Once Theorem 2.2 and an iterated in time discrete solution  $(\rho^n)_{n \geq 0}$  on hand, one can construct a piecewise constant in time and space reconstruction  $\rho_{\mathcal{T},\tau}$  by setting

$$(31) \quad \rho_{\mathcal{T},\tau}(t, x) = \rho_K^n \quad \text{if } (t, x) \in (t_{n-1}, t_n] \times K, \quad n \geq 1, \quad \rho_{\mathcal{T},\tau}(0, x) = \rho_K^0 \quad \text{if } x \in K.$$

Now, let  $(\mathcal{T}_m, \mathcal{E}_m, (x_K)_{K \in \mathcal{T}_m})_{m \geq 0}$  and  $(\tau_m)_{m \geq 0}$  be respectively a sequence of admissible meshes in the sense of Definition 2.1 and a sequence of time steps such that

$$(32) \quad \lim_{m \rightarrow \infty} \delta_{\mathcal{T}_m} = \lim_{m \rightarrow \infty} \tau_m = 0 \quad \text{and} \quad \zeta_{\mathcal{T}_m} \leq \zeta_\star < +\infty, \quad m \geq 0,$$

then the corresponding sequence of approximate solutions  $(\rho_{\mathcal{T}_m, \tau_m})_{m \geq 1}$  is bounded in  $L^\infty(\mathbb{R}_+ \times \Omega)$  owing to Theorem 2.2. Therefore, there exists  $\rho \in L^\infty(\mathbb{R}_+ \times \Omega)$  with  $0 \leq \rho \leq 1$  such that, up to a subsequence,

$$(33) \quad \rho_{\mathcal{T}_m, \tau_m} \xrightarrow{m \rightarrow \infty} \rho \quad \text{in the } L^\infty(\mathbb{R}_+ \times \Omega) \text{ weak-}\star \text{ sense.}$$

The following theorem claims that  $\rho$  is the unique weak solution to the continuous problem (1)–(4), and that the convergence holds in a stronger sense.

**Theorem 2.3.** *Let  $\rho$  be as in (33), then  $\rho$  is the unique weak solution to (1)–(4) in the sense of Definition 1.2. Moreover, the whole sequence  $(\rho_{\mathcal{T}_m, \tau_m})_{m \geq 0}$  converges strongly in  $L_{loc}^p(\mathbb{R}_+ \times \overline{\Omega})$  for any  $p \in [1, +\infty)$ .*

Proving Theorem 2.3 is the purpose of Section 4. The proof is based on compactness arguments that build on some refined version of the discrete energy estimate (30). Numerical evidences of the convergence will then be provided in Section 5.

### 3. NUMERICAL ANALYSIS AT FIXED GRID

The goal of this section is twofold. First one aims at establishing Theorem 2.2. Second, one derives enough estimates to carry out the convergence analysis in Section 4.

**3.1. Existence and uniqueness of the discrete solution.** We are interested in solutions  $\boldsymbol{\rho}^n$  to the scheme (18)–(20) that are bounded between 0 and 1. Therefore, changing the definition (19a) of the internal fluxes by

$$(34a) \quad F_{K\sigma}^n = \frac{1}{d_\sigma} \left[ (\rho_K^n)^+ (1 - \rho_L^n)^+ e^{\frac{1}{2}(\phi_K - \phi_L)} - (\rho_L^n)^+ (1 - \rho_K^n)^+ e^{\frac{1}{2}(\phi_L - \phi_K)} \right] \text{ for } \sigma = K|L \in \mathcal{E}_{\text{int}},$$

and the one (19b) of the boundary fluxes by

$$(34b) \quad F_{K\sigma}^n = \alpha_\sigma \rho_\sigma^n - \beta_\sigma \quad \text{with} \quad \rho_\sigma^n = \frac{d_\sigma \beta_\sigma + (\rho_K^n)^+ e^{\frac{1}{2}(\phi_K - \phi_\sigma)}}{d_\sigma \alpha_\sigma + (\rho_K^n)^+ e^{\frac{1}{2}(\phi_K - \phi_\sigma)} + (1 - \rho_K^n)^+ e^{-\frac{1}{2}(\phi_K - \phi_\sigma)}} \in (0, 1)$$

does not affect the value of the solution  $\boldsymbol{\rho}^n$ . After performing this slight modification, one can establish the following a priori estimate.

**Lemma 3.1.** *Given  $\boldsymbol{\rho}^{n-1} \in [0, 1]^\mathcal{T}$ ,  $n \geq 1$ , any solution  $\boldsymbol{\rho}^n$  to the modified scheme (18)&(34) belongs to  $(0, 1)^\mathcal{T}$ . In particular, being a solution to (18)&(34) is equivalent to being a solution in  $(0, 1)^\mathcal{T}$  to (18)–(20).*

*Proof.* We argue by contradiction. Assume that there exists  $K \in \mathcal{T}$  such that  $\rho_K^n \geq 1$ , then we deduce from (34) and from  $\alpha_\sigma > \beta_\sigma > 0$  that  $F_{K\sigma}^n \geq 0$  for all  $\sigma \in \mathcal{E}_K$ , and even that

$$(35) \quad F_{K\sigma}^n > 0 \quad \text{if} \quad \sigma \in \mathcal{E}_{\text{ext}}.$$

Then we deduce from (18) that

$$(36) \quad 0 \leq \sum_{\sigma \in \mathcal{E}_K} m_\sigma F_{K\sigma}^n = \frac{\rho_K^{n-1} - \rho_K^n}{\tau} m_K \leq 0.$$

Therefore,  $\rho_K^n = 1$  and all the fluxes  $F_{K\sigma}^n$ ,  $\sigma \in \mathcal{E}_K$  vanish. This is only possible if  $\rho_L^n = 1$  for each neighboring cell  $L$  such that  $\sigma = K|L \in \mathcal{E}_{\text{int}}$ . One can iterate to neighbors of neighbors until reaching  $K$  such that  $\mathcal{E}_{K,\text{ext}} \neq \emptyset$ . For such a cell, the first inequality in (36) is strict, leading to a contradiction, hence  $\rho_K^n < 1$  for all  $K \in \mathcal{T}$ . The bound  $\rho_K^n > 0$  for all  $K \in \mathcal{T}$  can be established in a similar way.  $\square$

**Proposition 3.2.** *For all  $n \geq 1$ , there exists a unique  $\boldsymbol{\rho}^n \in (0, 1)^\mathcal{T}$  solution to (16)–(19).*

*Proof.* The proof splits in two steps. Let us first show that at each time step  $n \geq 1$ , there exists a solution  $\boldsymbol{\rho}^n \in (0, 1)^\mathcal{T}$  to (18)&(19), or equivalently owing to Lemma 3.1,  $\boldsymbol{\rho}^n$  solution to (18)&(34). Note that here,  $\rho_\sigma^n$  is thought as a function of  $\rho_K^n$ , cf. (20), rather than as an independent unknown.

Let  $n \geq 1$  be such that  $\boldsymbol{\rho}^{n-1} \in [0, 1]^\mathcal{T}$  is given (this is the case for  $n = 1$  owing to (16)). For  $s \in [0, \tau]$ , define  $\boldsymbol{\rho}^{(s)} = \left( \rho_K^{(s)} \right)_{K \in \mathcal{T}}$  as a solution to

$$(37) \quad \left( \rho_K^{(s)} - \rho_K^{n-1} \right) m_K + s \sum_{\sigma \in \mathcal{E}_K} m_\sigma F_{K\sigma}^{(s)} = 0, \quad K \in \mathcal{T},$$

with  $F_{K\sigma}^{(s)}$  defined by (34) where  $\boldsymbol{\rho}^n$  has been replaced by  $\boldsymbol{\rho}^{(s)}$ . For  $s = 0$ , the above system rewrites  $\mathbb{M}\boldsymbol{\rho}^{(0)} = \mathbb{M}\boldsymbol{\rho}^{n-1}$ , with the matrix  $\mathbb{M} = \text{diag} \left( (m_K)_{K \in \mathcal{T}} \right)$  having a positive determinant. The unique solution  $\boldsymbol{\rho}^{(0)} = \boldsymbol{\rho}^{n-1}$  belongs to  $[0, 1]^\mathcal{T}$ , while any solution  $\boldsymbol{\rho}^{(s)}$  for  $s > 0$  belongs to  $(0, 1)^\mathcal{T}$  thanks to Lemma 3.1. A standard topological degree argument (see [20, 12] for a presentation of the topological gradient, and [13, 1] for applications in a similar context) then shows that the nonlinear system (37) admits at least one solution  $\boldsymbol{\rho}^{(s)}$  for all  $s > 0$ . In particular, for  $s = \tau$ , this shows the

existence of  $\boldsymbol{\rho}^n \in (0, 1)^{\mathcal{T}}$  solution (18)&(19). By a straightforward induction on  $n$ , one gets the existence of  $\boldsymbol{\rho}^n \in (0, 1)^{\mathcal{T}}$  for all  $n \geq 1$ .

The second step of the proof consists in proving uniqueness for the solution in  $(0, 1)^{\mathcal{T}}$  to (18)&(19). Since  $\mathbf{0} \leq \boldsymbol{\rho}^n \leq \mathbf{1}$ ,  $F_{K\sigma}^n$  is an increasing function of  $\rho_K^n$  and a non-increasing one of  $\rho_L^n$  for  $L \neq K$ . As a consequence, the nonlinear system corresponding to (18) rewrites

$$(38) \quad \mathcal{H}^n(\boldsymbol{\rho}^n) = \left( \mathcal{H}_K^n \left( \rho_K^n, (\rho_L^n)_{L \neq K} \right) \right)_{K \in \mathcal{T}} = \mathbf{0},$$

where  $\mathcal{H}_K$  is increasing w.r.t. its first variable and non-decreasing w.r.t. the others. Assume that the scheme admits another solution  $\check{\boldsymbol{\rho}}^n \in [0, 1]^{\mathcal{T}}$  corresponding to the same previous step data  $\boldsymbol{\rho}^{n-1}$ :

$$\mathcal{H}^n(\check{\boldsymbol{\rho}}^n) = \left( \mathcal{H}_K^n \left( \check{\rho}_K^n, (\check{\rho}_L^n)_{L \neq K} \right) \right)_{K \in \mathcal{T}} = \mathbf{0},$$

Therefore, denoting by  $a \wedge b = \min(a, b)$  and  $a \vee b = \max(a, b)$ , one has

$$\mathcal{H}_K^n \left( \rho_K^n, (\rho_L^n \wedge \check{\rho}_L^n)_{L \neq K} \right) \geq 0, \quad \mathcal{H}_K^n \left( \check{\rho}_K^n, (\rho_L^n \wedge \check{\rho}_L^n)_{L \neq K} \right) \geq 0,$$

so that, since  $\rho_K^n \wedge \check{\rho}_K^n$  is either equal to  $\rho_K^n$  or  $\check{\rho}_K^n$ ,

$$(39) \quad \mathcal{H}_K^n \left( \rho_K^n \wedge \check{\rho}_K^n, (\rho_L^n \wedge \check{\rho}_L^n)_{L \neq K} \right) \geq 0, \quad K \in \mathcal{T}.$$

Similarly, there holds

$$(40) \quad \mathcal{H}_K^n \left( \rho_K^n \vee \check{\rho}_K^n, (\rho_L^n \vee \check{\rho}_L^n)_{L \neq K} \right) \leq 0, \quad K \in \mathcal{T}.$$

Subtracting (39) to (40), summing over  $K \in \mathcal{T}$  and using the conservativity of the fluxes provides

$$\sum_{K \in \mathcal{T}} \frac{|\rho_K^n - \check{\rho}_K^n|}{\tau} m_K + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m_\sigma \alpha_\sigma |\rho_\sigma^n - \check{\rho}_\sigma^n| \leq 0,$$

where  $\check{\rho}_\sigma^n$  is computed thanks to formula (20) with  $\check{\rho}_K^n$  instead of  $\rho_K^n$ . We conclude that  $\boldsymbol{\rho}^n = \check{\boldsymbol{\rho}}^n$ , completing the proof of Proposition 3.2.  $\square$

**3.2. Discrete energy/dissipation estimates.** The goal of this section is to show some refined energy estimate implying in particular (30). We pay attention to precisely quantifying the energy dissipation since this information is key to derive the compactness results to be used in Section 4. Denote by  $\mathbf{F}^n = (F_{K\sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_K}$  the approximate fluxes at time step  $n \geq 1$ , then taking inspiration in [25, 26], we introduce the primal dissipation potential by setting

$$(41) \quad \mathcal{D}_{\mathcal{E}}(\boldsymbol{\rho}^n, \mathbf{F}^n) = \sum_{\sigma \in \mathcal{E}} a_\sigma \eta_\sigma^n \Psi \left( \frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n} \right) \geq 0,$$

where  $\Psi$  is the continuous nonnegative strictly convex even function vanishing at 0 with superlinear growth at  $\infty$  defined by

$$\Psi(z) = 2z \log \left( \frac{z + \sqrt{z^2 + 4}}{2} \right) - 2\sqrt{z^2 + 4} + 4, \quad z \in \mathbb{R},$$

and where  $\eta_\sigma^n$ , which is defined by (22) for  $\sigma \in \mathcal{E}_{\text{int}}$  and by  $\eta_\sigma^n = \sqrt{\eta(\rho_K^n)\eta(\rho_\sigma^n)}$  for  $\sigma \in \mathcal{E}_{\text{ext}}$ , is positive thanks to Proposition 3.2.

As highlighted by the notation, the dissipation is associated to the edges  $\mathcal{E}$ . Yet, the dissipation potential defined in (41) only corresponds to the dissipation in the bulk even though boundary fluxes also contribute to the dissipation of the total free energy, as shows (12). This choice is made

for simplicity and is possible since the quantification of the dissipation across the sole bulk already provides enough compactness to carry out the convergence proof, see Section 4. Note that each internal edge  $\sigma = K|L \in \mathcal{E}_{\text{int}}$  appears only once in (41) and that  $\Psi(\frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n}) = \Psi(\frac{d_\sigma F_{L\sigma}^n}{\eta_\sigma^n})$  since  $\Psi$  is even and since  $F_{K\sigma}^n + F_{L\sigma}^n = 0$ .

Given  $\mathbf{G}^n = (G_{K\sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_K}$  with  $G_{K\sigma}^n + G_{L\sigma}^n = 0$  for all  $\sigma = K|L \in \mathcal{E}_{\text{int}}$ , then we define the dual dissipation potential  $\mathcal{D}_\mathcal{E}^* : (0, 1)^\mathcal{T} \times \mathbb{R}^\mathcal{E} \rightarrow \mathbb{R}_+$  by

$$(42) \quad \mathcal{D}_\mathcal{E}^*(\boldsymbol{\rho}^n, \mathbf{G}^n) = \sum_{\sigma \in \mathcal{E}} a_\sigma \eta_\sigma^n \Psi^*(G_{K\sigma}^n) \geq 0,$$

where  $\Psi^*$  is the Legendre transform of  $\Psi$ , defined by

$$\Psi^*(s) = 4(\cosh(s/2) - 1), \quad s \in \mathbb{R}.$$

It is continuous, nonnegative, uniformly convex and vanishes at 0.

**Proposition 3.3.** *Let  $(\boldsymbol{\rho}^n)_{n \geq 1} \subset (0, 1)^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$  be the iterated solution to the scheme (16)–(20). For  $n \geq 1$ , let  $\mathbf{G}^n = (G_{K\sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_K}$  be defined by*

$$(43) \quad G_{K\sigma}^n = \xi_K^n - \xi_{K\sigma}^n = \begin{cases} \xi_K^n - \xi_L^n & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ \xi_K^n - \xi_\sigma^n & \text{if } \sigma \in \mathcal{E}_{\text{ext}}, \end{cases}$$

then there holds

$$(44) \quad \frac{\mathcal{F}_{\mathcal{T}, \text{tot}}^n - \mathcal{F}_{\mathcal{T}, \text{tot}}^{n-1}}{\tau} + \mathcal{D}_\mathcal{E}(\boldsymbol{\rho}^n, \mathbf{F}^n) + \mathcal{D}_\mathcal{E}^*(\boldsymbol{\rho}^n, \mathbf{G}^n) \leq 0.$$

*Proof.* Since the solution  $\boldsymbol{\rho}^n$ ,  $n \geq 1$ , belongs to  $(0, 1)^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$ , the discrete electrochemical potential  $\boldsymbol{\xi}^n = ((\xi_K^n)_{K \in \mathcal{T}}, (\xi_\sigma^n)_{\sigma \in \mathcal{E}_{\text{ext}}}) \in \mathbb{R}^{\mathcal{T} \cup \mathcal{E}_{\text{ext}}}$  is well defined. Multiplying the discrete conservation law (18) by  $\xi_K^n$  and summing over  $K \in \mathcal{T}$  leads to

$$(45) \quad A_{\mathcal{T}}^n + B_{\mathcal{T}}^n + C_{\mathcal{T}}^n = 0,$$

with

$$A_{\mathcal{T}}^n = \sum_{K \in \mathcal{T}} m_K \frac{\rho_K^n - \rho_K^{n-1}}{\tau} \xi_K^n, \quad B_{\mathcal{T}}^n = \sum_{\sigma \in \mathcal{E}} m_\sigma F_{K\sigma}^n G_{K\sigma}^n \quad \text{and} \quad C_{\mathcal{T}}^n = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m_\sigma F_{K\sigma}^n \xi_\sigma^n.$$

Similarly to what we did in (8) at the continuous level, the external fluxes  $F_{K\sigma}^n$  given by (19b) can be rewritten as

$$F_{K\sigma}^n = 2\sqrt{\beta_\sigma(\alpha_\sigma - \beta_\sigma)\rho_\sigma^n(1 - \rho_\sigma^n)} \sinh\left(\frac{1}{2}(\xi_\sigma^n - \xi_\sigma^\Gamma)\right), \quad \sigma \in \mathcal{E}_{\text{ext}}.$$

Therefore,  $F_{K\sigma}^n(\xi_\sigma^n - \xi_\sigma^\Gamma) \geq 0$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , so that

$$(46) \quad C_{\mathcal{T}}^n \geq \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m_\sigma F_{K\sigma}^n \xi_\sigma^\Gamma.$$

Concerning the bulk term  $B_{\mathcal{T}}^n$ , the writing (21) of the internal edge fluxes and its straightforward counterpart for boundary edges show that

$$\frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n} = 2 \sinh\left(\frac{G_{K\sigma}^n}{2}\right) = (\Psi^*)'(G_{K\sigma}^n).$$

Therefore, we have equality in the Young-Fenchel inequality

$$\frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n} G_{K\sigma}^n = \Psi \left( \frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n} \right) + \Psi^* (G_{K\sigma}^n).$$

As a consequence,

$$(47) \quad B_{\mathcal{T}}^n = \sum_{\sigma \in \mathcal{E}} a_\sigma \eta_\sigma^n \frac{d_\sigma F_{K\sigma}^n}{\eta_\sigma^n} G_{K\sigma}^n = \mathcal{D}_{\mathcal{E}}(\boldsymbol{\rho}^n, \mathbf{F}^n) + \mathcal{D}_{\mathcal{E}}^*(\boldsymbol{\rho}^n, \mathbf{G}^n).$$

For the accumulation term  $A_{\mathcal{T}}^n$ , the convexity of the mixing entropy density  $h$  implies that

$$(\rho_K^n - \rho_K^{n-1}) h'(\rho_K^n) \geq h(\rho_K^n) - h(\rho_K^{n-1}), \quad K \in \mathcal{T}.$$

Therefore, it follows from the definition (27) of  $\mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^n)$  that

$$(48) \quad A_{\mathcal{T}}^n \geq \frac{\mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^n) - \mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^{n-1})}{\tau}.$$

We recover the discrete energy dissipation estimate (44) by combining (46)–(48) in (45) and by using the definition (27) of  $\mathcal{F}_{\mathcal{T},\text{tot}}^n$ .  $\square$

Since  $\phi$  is assumed to be nonnegative, so does  $\mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^n)$ . The upper bound we deduce from Proposition 3.3 is rather on  $\mathcal{F}_{\mathcal{T},\text{tot}}^n$ , which is not bounded from below so far. Obtaining a time-dependent lower-bound for  $\mathcal{F}_{\mathcal{T},\text{tot}}^n$  is the purpose of the following corollary. Its proof, the details of which are left to the reader, relies on the fact that both  $\xi_\sigma^\Gamma$  and  $F_{K\sigma}^n$  are uniformly bounded for each  $\sigma \in \mathcal{E}_{\text{ext}}$ .

**Corollary 3.4.** *Let  $C_1$  be as in Proposition 1.1, then*

$$-C_1 t_n \leq \mathcal{F}_{\mathcal{T},\text{tot}}^n \leq \mathcal{F}_{\mathcal{T},\text{tot}}^{n-1} \leq \mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^0), \quad n \geq 1.$$

*In particular, (30) holds true.*

In Proposition 3.3, the free energy dissipation is quantified thanks to the non-homogeneous functionals  $\mathcal{D}_{\mathcal{E}}$  and  $\mathcal{D}_{\mathcal{E}}^*$ . The goal of the next Lemma is to deduce from this estimate some more classical discrete  $L_{\text{loc}}^2(\mathbb{R}_+; H^1(\Omega))$  estimate on  $(\boldsymbol{\rho}^n)_{n \geq 1}$ .

**Lemma 3.5.** *There exists  $C_2$  depending only on  $C_1$ ,  $\Omega$  and  $\phi$  such that*

$$\sum_{p=1}^n \tau \sum_{\sigma \in \mathcal{E}} a_\sigma (\rho_K^p - \rho_{K\sigma}^p)^2 \leq C_2(1 + t_n), \quad \forall n \geq 1.$$

*Proof.* Combining Proposition 3.3 with Corollary 3.4, we obtain that

$$(49) \quad \tau \sum_{p=1}^n \mathcal{D}_{\mathcal{E}}^*(\boldsymbol{\rho}^p, \mathbf{G}^p) = \sum_{p=1}^n \tau \sum_{\sigma \in \mathcal{E}} a_\sigma \eta_\sigma^p \Psi^*(G_{K\sigma}^p) \leq \mathcal{F}_{\mathcal{T},\text{tot}}^0 - \mathcal{F}_{\mathcal{T},\text{tot}}^n \leq \mathcal{F}_{\mathcal{T}}(\boldsymbol{\rho}^0) + C_1 t_n.$$

It follows from the elementary inequality  $\cosh(a+b) = \cosh(a)\cosh(b) + \sinh(a)\sinh(b)$  that

$$\begin{aligned} \eta_\sigma^p \Psi^*(G_{K\sigma}^p) &= 4\eta_\sigma^p \left( \cosh \left( \frac{\phi_K - \phi_{K\sigma}}{2} \right) \cosh \left( \frac{h'(\rho_K^p) - h'(\rho_{K\sigma}^p)}{2} \right) - 1 \right) \\ &\quad + 4\eta_\sigma^p \sinh \left( \frac{\phi_K - \phi_{K\sigma}}{2} \right) \sinh \left( \frac{h'(\rho_K^p) - h'(\rho_{K\sigma}^p)}{2} \right) =: S_\sigma^p + T_\sigma^p, \quad p \geq 1. \end{aligned}$$

Then using (24),  $\cosh(a) \geq 1$ , and the fact that the arithmetic mean is greater than the geometric one, one gets that

$$S_\sigma^p = 2 \cosh\left(\frac{\phi_K - \phi_{K\sigma}}{2}\right) ((\rho_K^p - \rho_{K\sigma}^p)^2 + \eta(\rho_K^p) + \eta(\rho_{K\sigma}^p) - 2\eta_\sigma^p) \geq 2(\rho_K^p - \rho_{K\sigma}^p)^2.$$

On the other hand, (23) yields

$$T_\sigma^p = 2(\rho_K^p - \rho_{K\sigma}^p) \sinh\left(\frac{\phi_K - \phi_{K\sigma}}{2}\right).$$

Since

$$\left|2 \sinh\left(\frac{\phi_K - \phi_{K\sigma}}{2}\right)\right| \leq \cosh\left(\frac{\|\phi\|_\infty}{2}\right) |\phi_K - \phi_{K\sigma}| \leq \cosh\left(\frac{\|\phi\|_\infty}{2}\right) \|\nabla\phi\|_\infty d_\sigma,$$

we deduce from Young's inequality that

$$T_\sigma^p \geq -(\rho_K^p - \rho_{K\sigma}^p)^2 - \cosh^2\left(\frac{\|\phi\|_\infty}{2}\right) \|\nabla\phi\|_\infty^2 d_\sigma^2.$$

All in all, we obtain

$$\begin{aligned} \tau \sum_{p=1}^n \mathcal{D}_\mathcal{E}^*(\rho^p, \mathbf{G}^p) &\geq \sum_{p=1}^n \tau \sum_{\sigma \in \mathcal{E}} a_\sigma (\rho_K^p - \rho_{K\sigma}^p)^2 - \cosh^2\left(\frac{\|\phi\|_\infty}{2}\right) \|\nabla\phi\|_\infty^2 \sum_{p=1}^n \tau \sum_{\sigma \in \mathcal{E}} m_\sigma d_\sigma \\ &= \sum_{p=1}^n \tau \sum_{\sigma \in \mathcal{E}} a_\sigma (\rho_K^p - \rho_{K\sigma}^p)^2 - \cosh^2\left(\frac{\|\phi\|_\infty}{2}\right) \|\nabla\phi\|_\infty^2 dm_\Omega t_n, \end{aligned}$$

which provides the desired result after being combined with (49).  $\square$

Next lemma exploits the other part of the dissipation to derive some discrete  $W_{\text{loc}}^{1,1}(\mathbb{R}_+; W^{-1,1}(\Omega))$  on  $(\rho^n)_{n \geq 0}$ .

**Lemma 3.6.** *Let  $\varphi \in C_c^\infty([0, T] \times \Omega)$  for some  $T > 0$ , with  $\text{dist}(\text{supp } \varphi, \partial\Omega) \geq \zeta_* \delta_\tau$ , then define  $\varphi_K^n = \frac{1}{m_K} \int_K \varphi(t_n)$  for all  $K \in \mathcal{T}$  and  $n \geq 0$ , then there exists  $C_3$  depending only on  $\Omega$ ,  $\alpha$ ,  $\beta$ ,  $\phi$ ,  $T$ , and  $\zeta_*$  such that*

$$\sum_{n \geq 1} \sum_{K \in \mathcal{T}} m_K (\rho_K^n - \rho_K^{n-1}) \varphi_K^n \leq C_3 \|\nabla\varphi\|_\infty.$$

*Proof.* The assumption on the support of  $\varphi$  implies that  $\varphi_K^n = 0$  either if  $K$  has a boundary edge  $\sigma \in \mathcal{E}_{K, \text{ext}}$  or if  $n \geq T/\tau$ . Moreover, it follows from the mean value theorem that for all  $K \in \mathcal{T}$  and all  $n \geq 1$ , there exists  $y_K^n \in K$  such that  $\varphi_K^n = \varphi(t_n, y_K^n)$ . Then for all  $\sigma = K|L \in \mathcal{E}_{\text{int}}$ , one has

$$|\varphi_K^n - \varphi_L^n| \leq \|\nabla\varphi\|_\infty (|y_K^n - x_K| + |y_L^n - x_L| + d_\sigma) \leq C_4 \|\nabla\varphi\|_\infty d_\sigma.$$

with  $C_4 = 1 + 2\zeta_\tau$ . Therefore, multiplying (18) by  $\tau\varphi_K^n$  and summing over  $K \in \mathcal{T}$  and  $n \geq 1$  provides

$$\begin{aligned} \sum_{n \geq 1} \sum_{K \in \mathcal{T}} m_K (\rho_K^n - \rho_K^{n-1}) \varphi_K^n &= - \sum_{n=1}^{\lfloor T/\tau \rfloor} \tau \sum_{\sigma \in \mathcal{E}_{\text{int}}} m_\sigma F_{K\sigma}^n (\varphi_K^n - \varphi_L^n) \\ &\leq C_4 \|\nabla\varphi\|_\infty \sum_{n=1}^{\lfloor T/\tau \rfloor} \tau \sum_{\sigma \in \mathcal{E}_{\text{int}}} a_\sigma \eta_\sigma^n \frac{d_\sigma |F_{K\sigma}^n|}{\eta_\sigma^n} d_\sigma, \end{aligned}$$

so that a Young-Fenchel inequality gives

$$(50) \quad \sum_{n \geq 1} \sum_{K \in \mathcal{T}} m_K (\rho_K^n - \rho_K^{n-1}) \varphi_K^n \leq C_4 \|\nabla \varphi\|_\infty \sum_{n=1}^{\lfloor T/\tau \rfloor} \tau \left( \mathcal{D}_\mathcal{E}(\boldsymbol{\rho}^n, \mathbf{F}^n) + \sum_{\sigma \in \mathcal{E}_{\text{int}}} a_\sigma \eta_\sigma^n \Psi^*(d_\sigma) \right).$$

A Taylor expansion of  $\Psi^*$  around 0 shows that

$$\Psi^*(d_\sigma) = \frac{d_\sigma^2}{2} (\Psi^*)''(c_\sigma) \quad \text{with } c_\sigma \in (0, d_\sigma) \subset [0, \text{diam}(\Omega)],$$

whence, since  $\eta_\sigma^n \leq 1/4$ ,

$$(51) \quad \sum_{\sigma \in \mathcal{E}_{\text{int}}} a_\sigma \eta_\sigma^n \Psi^*(d_\sigma) \leq \frac{1}{8} \cosh\left(\frac{\text{diam}(\Omega)}{2}\right) \sum_{\sigma \in \mathcal{E}_{\text{int}}} m_\sigma d_\sigma \leq \frac{d}{8} \cosh\left(\frac{\text{diam}(\Omega)}{2}\right) m_\Omega.$$

Then we deduce from Proposition 3.3 and Corollary 3.4 that

$$(52) \quad \sum_{n=1}^{\lfloor T/\tau \rfloor} \tau \mathcal{D}_\mathcal{E}(\boldsymbol{\rho}^n, \mathbf{F}^n) \leq \mathcal{F}_\mathcal{T}(\boldsymbol{\rho}^0) - \mathcal{F}_{\mathcal{T}, \text{tot}}^{\lfloor T/\tau \rfloor} \leq m_\Omega (\log 2 + \|\phi\|_\infty + 2\|\nabla \phi\|_\infty \delta_\mathcal{T}) + C_1 T.$$

The combination of (51)(52) in (50) shows the desired result.  $\square$

#### 4. CONVERGENCE ANALYSIS

The goal of this section is to prove Theorem 2.3. The proof consists in three steps. First in Section 4.1, we establish some compactness results on  $(\rho_{\mathcal{T}_m, \tau_m})_{m \geq 0}$ . Then we identify in Section 4.2 any limit value  $\rho$  of  $(\rho_{\mathcal{T}_m, \tau_m})_{m \geq 0}$  as a weak solution to the continuous problem. Finally, the uniqueness of the weak solution is established in Section 4.3, implying by the way the convergence of the whole sequence.

In what follows, we lighten the notation by removing the index  $m$  associated to the mesh and time step. The limit  $m \rightarrow +\infty$  is denoted by  $\delta_\mathcal{T}, \tau \rightarrow 0$  instead. This limit implicitly supposes that the regularity factor of the mesh  $\zeta_\mathcal{T}$  remains uniformly bounded by some  $\zeta_*$  as prescribed by (32).

**4.1. Compactness properties.** We derived in Section 3.2 all the preliminary material required to use some existing compactness results. First by combining Lemmas 3.5 and 3.6, one can apply the black-box discrete Aubin-Simon theorem [2, Theorem 3.9], leading to the following compactness result.

**Proposition 4.1.** *Let  $\rho$  be a limit value (33) of  $\rho_{\mathcal{T}, \tau}$  as  $\delta_\mathcal{T}, \tau$  tend to 0, then  $\rho \in L_{loc}^2(\mathbb{R}_+; H^1(\Omega))$  and, up to a subsequence,*

$$(53) \quad \rho_{\mathcal{T}, \tau} \xrightarrow{\delta_\mathcal{T}, \tau \rightarrow 0} \rho \quad \text{in } L_{loc}^p(\mathbb{R}_+ \times \Omega).$$

The above proposition shows some strong convergence in the bulk domain  $\mathbb{R}_+ \times \Omega$ . To pass to the limit in the boundary conditions, one also has to get some convergence of the traces on  $\mathbb{R}_+ \times \partial\Omega$ . Even though the boundary condition (3) in linear w.r.t.  $\rho$ , we establish the strong convergence of the trace of the approximate solution  $\rho_{\mathcal{T}, \tau}$  towards the trace of  $\rho$ .

**Lemma 4.2.** *Let  $\rho_{\mathcal{T}, \tau}$  be such that the convergence (53) holds. Denote by  $\gamma \rho_{\mathcal{T}, \tau}$  the trace on  $\mathbb{R}_+ \times \Gamma$  of the approximate solution  $\rho_{\mathcal{T}, \tau}$ , i.e.*

$$\gamma \rho_{\mathcal{T}, \tau}(t, x) = \rho_K^n \quad \text{for } (t, x) \in (t_{n-1}, t_n] \times \sigma, \quad \sigma \in \mathcal{E}_{K, \text{ext}}, \quad K \in \mathcal{T},$$

and by  $\gamma\rho \in L^2_{loc}(\mathbb{R}_+; H^{1/2}(\Gamma))$  the trace of a limit value  $\rho$  of  $\rho_{\mathcal{T},\tau}$ , then

$$(54) \quad \gamma\rho_{\mathcal{T},\tau} \xrightarrow{\delta_{\mathcal{T},\tau} \rightarrow 0} \gamma\rho \quad \text{in } L^p_{loc}(\mathbb{R}_+ \times \Gamma), \quad 1 \leq p < +\infty.$$

*Proof.* The proof builds on ideas introduced in Section 4.2 of [4]. First, notice that since both  $\rho_{\mathcal{T},\tau}$  and  $\rho$  remain bounded between 0 and 1, it suffices to establish the convergence (54) in  $L^2_{loc}(\mathbb{R}_+ \times \Gamma)$  to get it all the  $L^p_{loc}(\mathbb{R}_+ \times \Gamma)$  thanks to the dominated convergence theorem.

Since  $\Omega$  is assumed to be polyhedral, its boundary  $\Gamma$  can be decomposed as  $\Gamma = \bigcup_{i=1}^I \Gamma_i$  with  $\Gamma_i$  included in an hyperplane of  $\mathbb{R}^d$  and  $I$  finite. We assume that the  $\Gamma_i$  are disjointed one from another. For  $\varepsilon > 0$  and  $i \in \{1, \dots, I\}$ , we define

$$\Gamma_{i,\varepsilon} = \{x \in \Gamma_i \mid x - \theta\nu_i \in \Omega \text{ for } \theta \in [0, \varepsilon]\},$$

where  $\nu_i$  is the outward w.r.t.  $\Omega$  normal to  $\Gamma_i$ . Denoting by  $m_{\Gamma_{i,\varepsilon}}$  (resp.  $m_{\Gamma_i}$ ) the  $(d-1)$ -dimensional Hausdorff (or Lebesgue) measure of  $\Gamma_{i,\varepsilon}$  (resp.  $\Gamma_i$ ), then

$$m_{\Gamma_i} - C_5\varepsilon \leq m_{\Gamma_{i,\varepsilon}} \leq m_{\Gamma_i}, \quad \varepsilon > 0, \quad 1 \leq i \leq I,$$

for some  $C_5$  depending only on  $\Omega$ . Therefore, given an arbitrary final time  $T > 0$  and an arbitrary  $\varepsilon > 0$ , then for all  $i \in \{1, \dots, I\}$ , there holds

$$(55) \quad \int_0^T \int_{\Gamma_i} |\gamma\rho_{\mathcal{T},\tau} - \gamma\rho|^2 \leq \int_0^T \int_{\Gamma_{i,\varepsilon}} |\gamma\rho_{\mathcal{T},\tau} - \gamma\rho|^2 + C_5\varepsilon T.$$

Using  $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$ , we obtain that

$$(56) \quad \int_0^T \int_{\Gamma_{i,\varepsilon}} |\gamma\rho_{\mathcal{T},\tau}(t, y) - \gamma\rho(t, y)|^2 dy dt \leq A_{\mathcal{T},\tau}^\varepsilon + B_{\mathcal{T},\tau}^\varepsilon + C^\varepsilon,$$

with

$$\begin{aligned} A_{\mathcal{T},\tau}^\varepsilon &= \frac{3}{\varepsilon} \int_0^T \int_0^\varepsilon \int_{\Gamma_{i,\varepsilon}} |\gamma\rho_{\mathcal{T},\tau}(t, y) - \rho_{\mathcal{T},\tau}(t, y - \theta\nu_i)|^2 dy d\theta dt, \\ B_{\mathcal{T},\tau}^\varepsilon &= \frac{3}{\varepsilon} \int_0^T \int_0^\varepsilon \int_{\Gamma_{i,\varepsilon}} |\rho_{\mathcal{T},\tau}(t, y - \theta\nu_i) - \rho(t, y - \theta\nu_i)|^2 dy d\theta dt, \\ C^\varepsilon &= \frac{3}{\varepsilon} \int_0^T \int_0^\varepsilon \int_{\Gamma_{i,\varepsilon}} |\gamma\rho(t, y) - \rho(t, y - \theta\nu_i)|^2 dy d\theta dt. \end{aligned}$$

First, applying Lemma 4.8 of [4] in combination with Lemma 3.5 yields

$$(57) \quad A_{\mathcal{T},\tau}^\varepsilon \leq 3(\varepsilon + \delta_{\mathcal{T}}) \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{\sigma \in \mathcal{E}_{\text{int}}} a_\sigma (\rho_K^n - \rho_L^n)^2 \leq 3C_2(1 + T + \tau)(\varepsilon + \delta_{\mathcal{T}}).$$

Second, it results from Proposition 4.1 that, for any fixed  $\varepsilon > 0$ , there holds

$$(58) \quad \lim_{\delta_{\mathcal{T},\tau} \rightarrow 0} B_{\mathcal{T},\tau}^\varepsilon = 0.$$

Putting (55)–(58) altogether leads to

$$(59) \quad \limsup_{\delta_{\mathcal{T},\tau} \rightarrow 0} \int_0^T \int_{\Gamma_i} |\gamma\rho_{\mathcal{T},\tau} - \gamma\rho|^2 \leq (C_5T + 3C_2(1 + T))\varepsilon + C^\varepsilon, \quad \forall \varepsilon > 0.$$



Eventually, one lets  $\varepsilon \rightarrow 0$  in (59), the right-hand side of which and in particular  $C^\varepsilon$  tend to 0 since  $\gamma\rho$  is the trace of  $\rho$ . This concludes the proof of Lemma 4.2.  $\square$

Even though the term trace is slightly abusive, it is natural to introduce the alternative notion of trace on  $\mathbb{R}_+ \times \Gamma$  for the approximate solution  $\rho_{\mathcal{T},\tau}$  by setting

$$\tilde{\gamma}\rho_{\mathcal{T},\tau}(t,x) = \rho_\sigma^n \quad \text{for } (t,x) \in (t_{n-1}, t_n] \times \sigma, \quad \sigma \in \mathcal{E}_{\text{ext}}.$$

**Lemma 4.3.** *Let  $\rho_{\mathcal{T},\tau}$  be such that the convergence (53) holds, then for all  $T > 0$ , there holds*

$$(60) \quad \|\gamma\rho_{\mathcal{T},\tau} - \tilde{\gamma}\rho_{\mathcal{T},\tau}\|_{L^p((0,T)\times\Gamma)} \xrightarrow{\delta_{\mathcal{T},\tau} \rightarrow 0} 0, \quad 1 \leq p < +\infty.$$

In particular,  $\tilde{\gamma}\rho_{\mathcal{T},\tau}$  also tends to  $\gamma\rho$  in  $L^p_{\text{loc}}(\mathbb{R}_+ \times \Gamma)$  for all finite  $p$ .

*Proof.* Once again, the uniform  $L^\infty$  bounds on  $\gamma\rho_{\mathcal{T},\tau}$  and  $\tilde{\gamma}\rho_{\mathcal{T},\tau}$  allow to establish (60) for  $p = 1$  only. Then, going back to the definitions of  $\gamma\rho_{\mathcal{T},\tau}$  and  $\tilde{\gamma}\rho_{\mathcal{T},\tau}$ , Cauchy-Schwarz inequality gives

$$\begin{aligned} \|\gamma\rho_{\mathcal{T},\tau} - \tilde{\gamma}\rho_{\mathcal{T},\tau}\|_{L^1((0,T)\times\Gamma)}^2 &\leq \left( \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma |\rho_K^n - \rho_\sigma^n| \right)^2 \\ &\leq \left( \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} a_\sigma (\rho_K^n - \rho_\sigma^n)^2 \right) \left( \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma d_\sigma \right). \end{aligned}$$

Thanks to Lemma 3.5, the first term of the right-hand side can be overestimated by

$$\sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} a_\sigma (\rho_K^n - \rho_\sigma^n)^2 \leq \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{\sigma \in \mathcal{E}} a_\sigma (\rho_K^n - \rho_{K\sigma}^n)^2 \leq C_2(1 + T + \tau).$$

On the other hand, it follows from the regularity of the mesh that

$$\sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma d_\sigma \leq \zeta_* \delta_{\mathcal{T}} \sum_{n=1}^{\lceil T/\tau \rceil} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma \leq \zeta_* \delta_{\mathcal{T}} (T + \tau) m_\Gamma$$

where  $m_\Gamma$  denote the  $(d-1)$ -dimensional Hausdorff (or Lebesgue) measure of  $\Gamma$ . In particular, (60) holds for  $p = 1$ , and thus also for all finite  $p$ . The last statement of the lemma, namely the convergence of  $\tilde{\gamma}\rho_{\mathcal{T},\tau}$  towards  $\gamma\rho$ , is then a straightforward consequence of Lemma 4.2.  $\square$

**4.2. Identification of the limit.** Our goal is here to establish the consistency of the scheme by identifying any limit value  $\rho$  of  $\rho_{\mathcal{T},\tau}$  as a solution to the continuous problem.

**Proposition 4.4.** *Let  $\rho$  be a limit value of  $\rho_{\mathcal{T},\tau}$  as  $\delta_{\mathcal{T}}, \tau$  tend to 0, then  $\rho$  is a weak solution to the problem (1)–(4) in the sense of Definition 1.2.*

*Proof.* Let  $\varphi \in C_c^\infty(\mathbb{R}_+ \times \bar{\Omega})$ , then define  $\varphi_K^n = \varphi(x_K, t_n)$  and  $\varphi_\sigma^n = \varphi(x_\sigma, t_n)$  for all  $K \in \mathcal{T}$ , all  $\sigma \in \mathcal{E}_{\text{ext}}$  and  $n \geq 0$ . This allows to define the function  $\varphi_{\mathcal{T},\tau}$  by

$$\varphi_{\mathcal{T},\tau}(t,x) = \varphi_K^{n-1} \quad \text{if } (t,x) \in [t_{n-1}, t_n) \times K.$$

Multiplying (18) by  $\tau\varphi_K^{n-1}$  and summing over  $K \in \mathcal{T}$  provides

$$(61) \quad A_{\mathcal{T},\tau} + B_{\mathcal{T},\tau} = 0,$$

where we have set

$$A_{\mathcal{T},\tau} = \sum_{n \geq 1} \sum_{K \in \mathcal{T}} m_K (\rho_K^n - \rho_K^{n-1}) \varphi_K^{n-1}, \quad B_{\mathcal{T},\tau} = \sum_{n \geq 1} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma F_{K\sigma}^n \varphi_K^{n-1}.$$

Since  $\varphi_K^n = 0$  for  $n$  large enough, the term  $A_{\mathcal{T},\tau}$  can be rewritten as

$$A_{\mathcal{T},\tau} = - \sum_{n \geq 1} \tau \sum_{K \in \mathcal{T}} m_K \rho_K^n \frac{\varphi_K^n - \varphi_K^{n-1}}{\tau} - \sum_{K \in \mathcal{T}} m_K \rho_K^0 \varphi_K^0.$$

Then classical arguments (see for instance [15]) allow to show that

$$(62) \quad \lim_{\delta_{\mathcal{T},\tau} \rightarrow 0} A_{\mathcal{T},\tau} = - \iint_{\mathbb{R}_+ \times \Omega} \rho \partial_t \varphi - \int_{\Omega} \rho^0 \varphi(0).$$

On the other hand, thanks to the conservativity of the fluxes, the term  $B_{\mathcal{T},\tau}$  reformulates as

$$B_{\mathcal{T},\tau} = \sum_{n \geq 1} \tau \sum_{\sigma \in \mathcal{E}} m_\sigma F_{K\sigma}^n (\varphi_K^{n-1} - \varphi_{K\sigma}^{n-1}) + \sum_{n \geq 1} \tau \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma F_{K\sigma}^n \varphi_\sigma^{n-1} =: B_{\mathcal{T},\tau}^{\text{bulk}} + B_{\mathcal{T},\tau}^{\text{ext}}.$$

Using the expression of the boundary fluxes (19b) in the term  $B_{\mathcal{T},\tau}^{\text{ext}}$  provides

$$B_{\mathcal{T},\tau}^{\text{ext}} = \iint_{\mathbb{R}_+ \times \Gamma} (\alpha_{\mathcal{E}} \tilde{\gamma} \rho_{\mathcal{T},\tau} - \beta_{\mathcal{E}}) \tilde{\gamma} \varphi_{\mathcal{T},\tau}$$

where  $\alpha_{\mathcal{E}}$  and  $\beta_{\mathcal{E}}$  are the piecewise constant (per edges  $\sigma \in \mathcal{E}_{\text{ext}}$ ) reconstructions on  $\Gamma$  build from the evaluation of  $\alpha$  and  $\beta$  at  $x_\sigma$ , and where

$$\tilde{\gamma} \varphi_{\mathcal{T},\tau}(t, x) = \varphi_\sigma^{n-1} \quad \text{if } (t, x) \in [t_{n-1}, t_n] \times \sigma, \quad \sigma \in \mathcal{E}_{\text{ext}}.$$

Due to the Lipschitz regularity of  $\alpha, \beta$  and  $\varphi$ , their approximations  $\alpha_{\mathcal{E}}, \beta_{\mathcal{E}}$  and  $\tilde{\gamma} \varphi_{\mathcal{T},\tau}$  converge uniformly. One concludes from the convergence of  $\tilde{\gamma} \rho_{\mathcal{T},\tau}$  stated at Lemma 4.3 that

$$(63) \quad \lim_{\delta_{\mathcal{T},\tau} \rightarrow 0} B_{\mathcal{T},\tau}^{\text{ext}} = \iint_{\mathbb{R}_+ \times \Gamma} (\alpha \gamma \rho - \beta) \varphi.$$

For the term  $B_{\mathcal{T},\tau}^{\text{bulk}}$ , we use the expression (25) of the internal fluxes, leading to

$$(64) \quad B_{\mathcal{T},\tau}^{\text{bulk}} = B_{\mathcal{T},\tau}^{\text{diff}} + B_{\mathcal{T},\tau}^{\text{conv}} + R_{\mathcal{T},\tau},$$

with

$$\begin{aligned} B_{\mathcal{T},\tau}^{\text{diff}} &= \sum_{n \geq 1} \tau \sum_{\sigma \in \mathcal{E}} a_\sigma (\rho_K^n - \rho_{K\sigma}^n) (\varphi_K^{n-1} - \varphi_{K\sigma}^{n-1}), \\ B_{\mathcal{T},\tau}^{\text{conv}} &= \sum_{n \geq 1} \tau \sum_{\sigma \in \mathcal{E}} a_\sigma \frac{\eta(\rho_K^n) + \eta(\rho_{K\sigma}^n)}{2} (\phi_K - \phi_{K\sigma}) (\varphi_K^{n-1} - \varphi_{K\sigma}^{n-1}), \\ R_{\mathcal{T},\tau} &= \sum_{n \geq 1} \tau \sum_{\sigma \in \mathcal{E}} m_\sigma R_\sigma^n (\varphi_K^{n-1} - \varphi_{K\sigma}^{n-1}). \end{aligned}$$

We do not detail the proof of

$$(65) \quad B_{\mathcal{T},\tau}^{\text{diff}} \xrightarrow{\delta_{\mathcal{T},\tau} \rightarrow 0} \iint_{\mathbb{R}_+ \times \Omega} \nabla \rho \cdot \nabla \varphi, \quad B_{\mathcal{T},\tau}^{\text{conv}} \xrightarrow{\delta_{\mathcal{T},\tau} \rightarrow 0} \iint_{\mathbb{R}_+ \times \Omega} \eta(\rho) \nabla \phi \cdot \nabla \varphi,$$

since similar terms have been studied in many contributions, see for instance [8] and references therein. It remains to show that  $R_{\mathcal{T},\tau}$  vanishes at the limit. We deduce from the expression (26), from the fact that  $r_\sigma^n \in (0, 1)$ , and from  $\|\eta\|_\infty = 1/4$  that

$$|R_\sigma^n| \leq \frac{1}{2} \|\nabla\phi\|_\infty (\rho_K^n - \rho_{K\sigma}^n)^2 + \frac{d_\sigma}{8} \|\nabla\phi\|_\infty^2 \left( |\rho_K^n - \rho_{K\sigma}^n| \cosh\|\phi\|_\infty + d_\sigma \|\nabla\phi\|_\infty \left( \frac{1}{2} + (\rho_K^n - \rho_{K\sigma}^n)^2 \right) \frac{1}{2} \cosh\|\phi\|_\infty \right).$$

Therefore, using furthermore Lemma 3.5, one readily shows that

$$(66) \quad R_{\mathcal{T},\tau} \leq C\delta_{\mathcal{T}}^2 \xrightarrow{\delta_{\mathcal{T},\tau} \rightarrow 0} 0.$$

Putting (62)–(66) together in (61) concludes the proof of Proposition 4.4.  $\square$

**4.3. Uniqueness of the weak solution.** So far, we established the convergence of the scheme towards a weak solution up to a subsequence. In order to show that the whole sequence converges, it suffices to show that the limit value is unique. This is a consequence of the following proposition.

**Proposition 4.5.** *The weak solution  $\rho$  to (1)–(4) in the sense of Definition 1.2 is unique.*

*Proof.* Let  $\rho$  and  $\check{\rho}$  be two weak solutions corresponding to the same initial data  $\rho^0$ , and let  $T$  be an arbitrary time horizon, then subtracting their respective weak formulations leads to

$$(67) \quad \int_0^T \langle \partial_t(\rho - \check{\rho}), \varphi \rangle_{H^{-1}, H_0^1} + \int_0^T \int_\Omega ((\eta(\rho) - \eta(\check{\rho})) \nabla\phi + \nabla(\rho - \check{\rho})) \cdot \nabla\varphi = 0$$

for all  $\varphi \in L^2((0, T); H_0^1(\Omega))$ . Choose  $\varphi$  as the solution to

$$-\Delta\varphi(t) = \rho(t) - \check{\rho}(t) \quad \text{in } \Omega, \quad \varphi(t) = 0 \quad \text{on } \Gamma, \quad t \in [0, T],$$

then one readily checks that  $\|\nabla\varphi(t)\|_{L^2(\Omega)^d} = \|\rho(t) - \check{\rho}(t)\|_{H^{-1}(\Omega)}$ . Moreover,  $\partial_t\varphi$  also belongs to  $L^2((0, T); H_0^1(\Omega))$  since  $\partial_t(\rho - \check{\rho})$  belongs to  $L^2((0, T); H^{-1}(\Omega))$ . Therefore,

$$\int_0^T \langle \partial_t(\rho - \check{\rho}), \varphi \rangle_{H^{-1}, H_0^1} = \int_0^T \int_\Omega \partial_t \nabla\varphi \cdot \nabla\varphi = \frac{1}{2} \|\nabla\varphi(T)\|_{(L^2(\Omega))^d}^2$$

since  $\varphi(0) = 0$ . As a consequence, (67) yields

$$\begin{aligned} \frac{1}{2} \|\rho(T) - \check{\rho}(T)\|_{H^{-1}(\Omega)}^2 + \|\rho - \check{\rho}\|_{L^2((0, T) \times \Omega)}^2 &= - \int_0^T \int_\Omega (\eta(\rho) - \eta(\check{\rho})) \nabla\phi \cdot \nabla\varphi \\ &\leq \|\nabla\phi\|_\infty \int_0^T \int_\Omega |\rho - \check{\rho}| |\nabla\varphi|. \end{aligned}$$

Then we deduce from Young's inequality that

$$\|\rho(T) - \check{\rho}(T)\|_{H^{-1}(\Omega)}^2 \leq \frac{\|\nabla\phi\|_\infty^2}{2} \|\rho - \check{\rho}\|_{L^2((0, T); H^{-1}(\Omega))}^2.$$

The above inequality holds for all  $T \geq 0$ , and we deduce from Gronwall Lemma together and from the fact that  $\rho(0) = \check{\rho}(0) = \rho^0$  that  $\|\rho(T) - \check{\rho}(T)\|_{H^{-1}(\Omega)} = 0$  for all  $T \geq 0$ .  $\square$

## 5. NUMERICAL RESULTS

Before presenting numerical results, let us comment briefly on some practical details concerning the effective implementation. Our code is based on Matlab. The resolution of the nonlinear system (18)–(19), which rewrites under the compact form (38) is achieved thanks to Newton's method:

$$(68) \quad \mathbb{J}(\boldsymbol{\rho}^{n,\ell})\delta\boldsymbol{\rho}^{n,\ell} = -\mathcal{H}^n(\boldsymbol{\rho}^{n,\ell}), \quad \boldsymbol{\rho}^{n,\ell+1} = \boldsymbol{\rho}^{n,\ell} + \delta\boldsymbol{\rho}^{n,\ell},$$

with  $\mathbb{J}$  standing for the Jacobian matrix of  $\mathcal{H}^n$ . Note that  $\rho_\sigma^n$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$  is not considered as an unknown and is deduced from the cell values thanks to (20). We initialize (68) by setting  $\boldsymbol{\rho}^{n,0} = \boldsymbol{\rho}^{n-1}$  and then iterate until  $\|\delta\boldsymbol{\rho}^{n,\ell}\|_\infty / \|\boldsymbol{\rho}^{n,\ell+1}\|_\infty \leq 10^{-12}$ . Then we set  $\boldsymbol{\rho}^n = \boldsymbol{\rho}^{n,\ell+1}$ .

**5.1. Numerical evidence of the convergence.** The first numerical test we propose aims at confirming our intuition concerning the second order accuracy in space of the scheme sketched in Section 2.2. To this end, we consider a one-dimensional domain  $\Omega = (0, 1)$ . We consider a slightly more general case than the one addressed in the paper by introducing some parameter  $\epsilon > 0$  (referred later on as the inverse Péclet number) in front of the diffusion term in (1b):

$$(69) \quad F + \eta(\rho)\partial_x\phi + \epsilon\partial_x\rho = 0.$$

The bulk numerical flux formula (19a) is tuned into

$$(70) \quad F_{K\sigma}^n = \frac{\epsilon}{d_\sigma} \left[ \rho_K^n (1 - \rho_L^n) e^{\frac{1}{2\epsilon}(\phi_K - \phi_L)} - \rho_L^n (1 - \rho_K^n) e^{\frac{1}{2\epsilon}(\phi_L - \phi_K)} \right] \quad \text{for } \sigma = K|L \in \mathcal{E}_{\text{int}}.$$

The boundary condition (3) remains unchanged at the continuous level, yet the discrete external fluxes are modified into

$$(71) \quad F_{K\sigma}^n = \frac{\epsilon}{d_\sigma} \left[ \rho_K^n (1 - \rho_\sigma^n) e^{\frac{1}{2\epsilon}(\phi_K - \phi_\sigma)} - \rho_\sigma^n (1 - \rho_K^n) e^{\frac{1}{2\epsilon}(\phi_\sigma - \phi_K)} \right] = \alpha_\sigma \rho_\sigma^n - \beta_\sigma, \quad \text{for } \sigma \in \mathcal{E}_{\text{int}},$$

with the updated boundary density value

$$(72) \quad \rho_\sigma^n = \frac{d_\sigma \beta_\sigma + \epsilon \rho_K^n e^{\frac{1}{2\epsilon}(\phi_K - \phi_\sigma)}}{d_\sigma \alpha_\sigma + \epsilon \rho_K^n e^{\frac{1}{2\epsilon}(\phi_K - \phi_\sigma)} + \epsilon (1 - \rho_K^n) e^{-\frac{1}{2\epsilon}(\phi_K - \phi_\sigma)}}.$$

The extension of our analysis to this framework is straightforward for fixed values of  $\epsilon > 0$ . In our test case, the functions  $\alpha$  and  $\beta$  defined on  $\Gamma = \{0, 1\}$  are chosen constant, with  $\alpha = 1$  and  $\beta = 1/2$ . Concerning the external potential, we set  $\phi(x) = 1 - x$ , so that the drift  $\partial_x\phi$  is constant. As an initial data, we choose

$$\rho^0(x) = \begin{cases} 1 & \text{if } x < 1/2, \\ 0 & \text{otherwise.} \end{cases}$$

The domain  $\Omega$  is discretized with a successively refined uniform grid. The final time is set to  $T = 2$ , whereas the time step  $\tau = 10^{-2}$  remains unchanged, in opposition to the spatial mesh size. A reference solution is computed on a fine grid made of 51200 cells.

We illustrate on Figure 1 the second order convergence in space that was expected from the discussion of Section 2.2. One notices that the error increases when the inverse Péclet number decreases. To better illustrate this point, we plot on Figure 2 the evolution of the error as a function of  $\epsilon$ . Such a behavior is expected since the scheme is not asymptotic preserving in the sense that the scheme corresponding to the limit  $\epsilon = 0$  is not consistent with the limiting hyperbolic continuous equation.

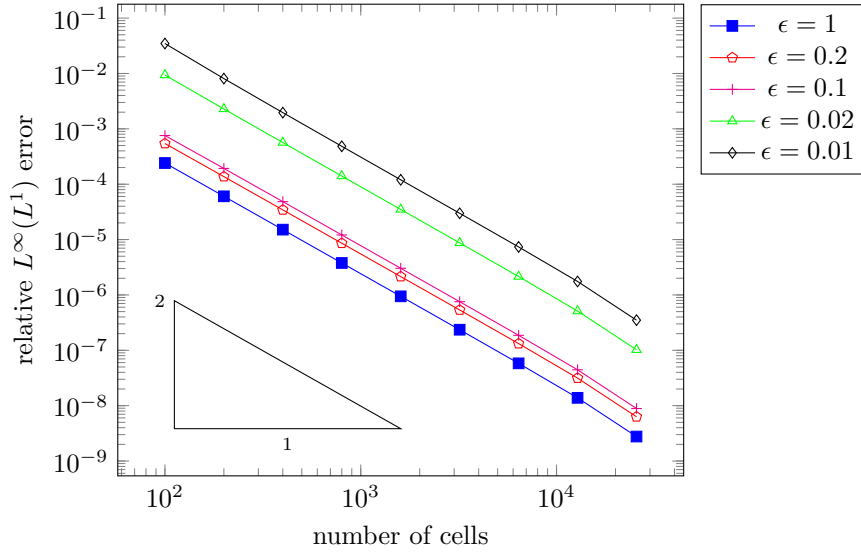


FIGURE 1. Evolution  $L^\infty((0, T); L^1(\Omega))$  relative errors as a function of the number of cells in the spatial discretization for various inverse Péclet numbers  $\epsilon$ .

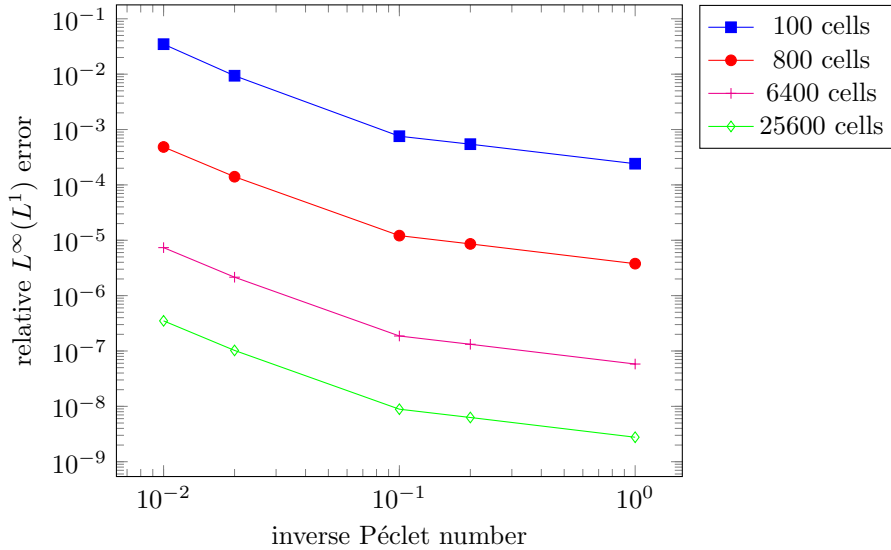


FIGURE 2. Evolution  $L^\infty((0, T); L^1(\Omega))$  on different meshes depending on the inverse Péclet number  $\epsilon$ .

**5.2. Energy stability and long-time behavior.** Our second numerical experiment is performed on a 2D Delaunay mesh made of 7374 triangles. Our goal is here twofold. First, we give a numerical evidence of the fact that the total energy  $\mathcal{F}_{\text{tot}}$  decreases along time, while the bulk energy  $\mathcal{F}(\rho)$

remains bounded. As in Section 5.1, we introduce the inverse Péclet number, which is set to  $\epsilon = 0.01$  in what follows. Therefore, the energy has to be adapted accordingly by setting

$$\mathcal{F}(\rho) = \int_{\Omega} (\epsilon h(\rho) + \rho \phi)$$

with  $\phi(x) = 1 - x_2$  for  $x = (x_1, x_2) \in \Omega$ . Concerning the boundary conditions,  $\alpha = 1$  is set constant, whereas

$$\beta(x) = \frac{1}{10} + \frac{4}{5} \left( \cos^2 \left( \frac{3\pi x_2}{2} \right) + (2x_2 - 1) \sin(\pi x_1) \right), \quad x = (x_1, x_2) \in \Omega.$$

As an initial data, we choose  $\rho^0(x) = 1$  if  $x \in (0, 1/2) \times (0, 1/2)$  and  $\rho^0(x) = 0$  otherwise. Snapshots of the solution are presented on Figure 3.

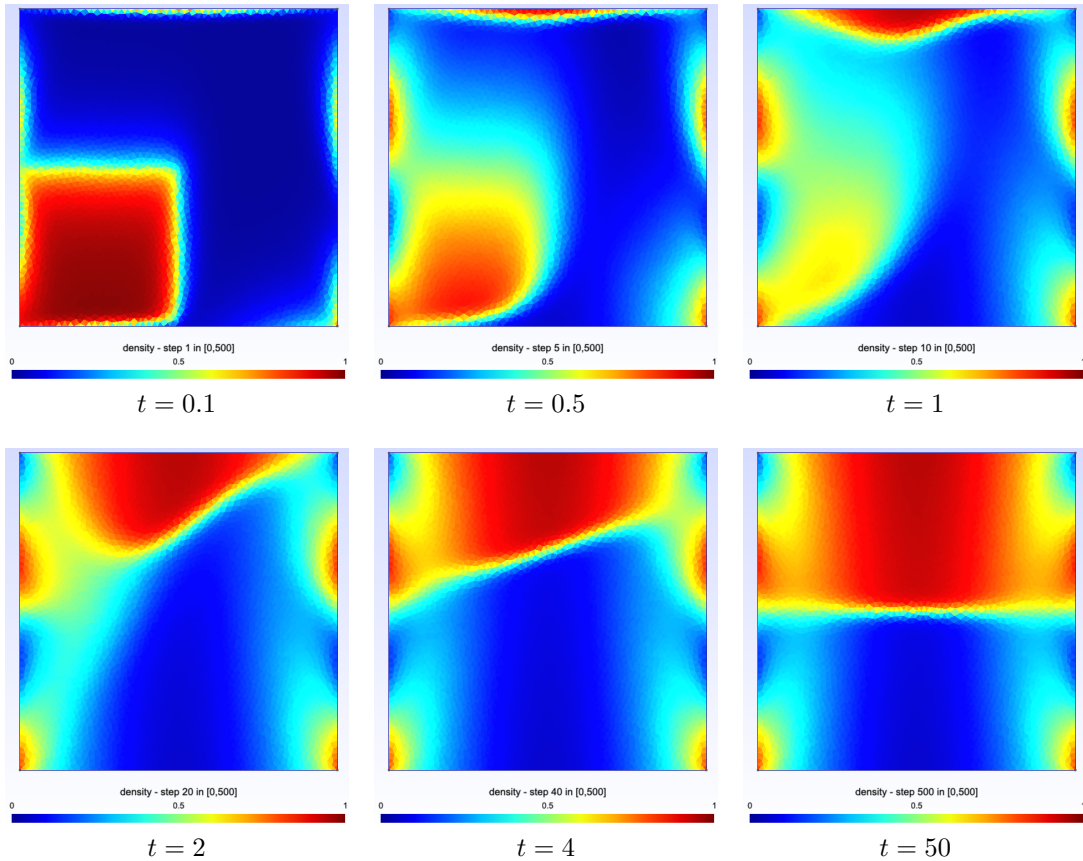


FIGURE 3. Snapshots of the solution at different times.

We plot on Figure 4 the evolution of the bulk and total energies along time. As expected,  $\mathcal{F}_{\text{tot}}$  is decreasing with linear decay, while  $\mathcal{F}(\rho)$  remains bounded along time.

We make use of a uniform time step  $\tau = 0.1$  until we reach the final time  $T = 50$ .

Second, we highlight the good behavior of the numerical scheme when it comes to the effective resolution of the induced nonlinear system. As expected, the highest number of required Newton

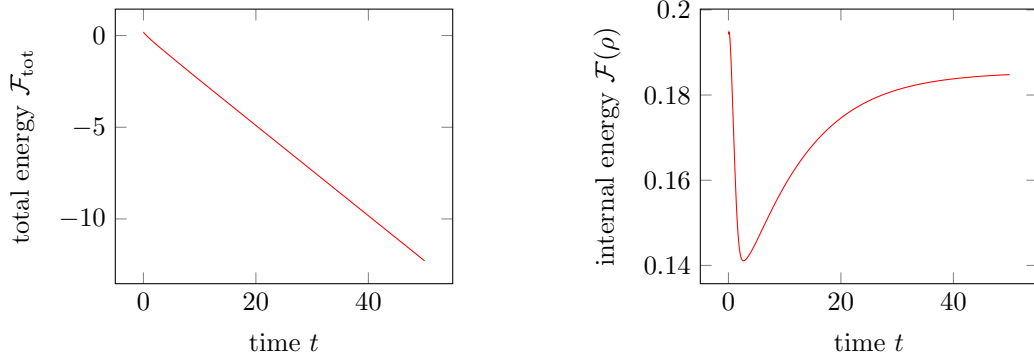


FIGURE 4. Evolution of the total energy (left) and of the bulk energy (right) along time.

iteration corresponds to the initial time steps where only 17 Newton iterations are required although  $\rho^1$  significantly differs from  $\rho^0$ . As time goes, this number decreases. In our test case, the steady state is not yet reached for  $T = 50$  and still 9 Newton iterations per time step are needed to solve the nonlinear system. This number can be importantly decreased for less demanding stopping criteria. The number of required Newton iterations at each time step is reported on Figure 5.

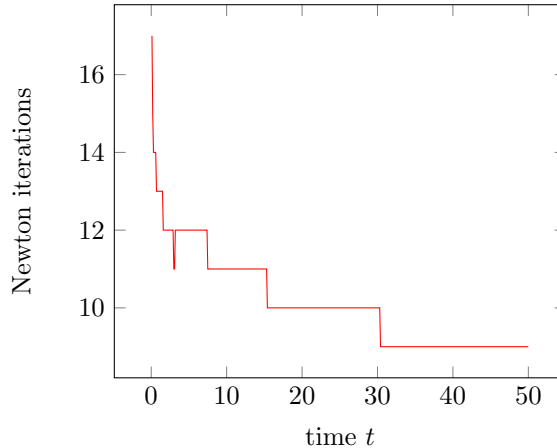


FIGURE 5. Number of Newton iterations required to solve the nonlinear system at each time step.

**Acknowledgements.** This project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 847593 (EURAD program, WP DONUT), and was further supported by Labex CEMPI (ANR-11-LABX-0007-01) and the Fédération de Recherche Mathématique des Hauts-de-France (FR CNRS 2037). C. Cancès also acknowledges support from the COMODO (ANR-19-CE46-0002) and MICMOV (ANR-19-CE40-0012) projects. J. Venel warmly thanks the Inria research center of the University of Lille for its hospitality and support, and the authors further thank Claire Chainais for stimulating discussions.

## REFERENCES

- [1] A. Ait Hammou Oulhaj. Numerical analysis of a finite volume scheme for a seawater intrusion model with cross-diffusion in an unconfined aquifer. *Numer. Methods Partial Differential Equations*, 34(3):857–880, 2018.
- [2] B. Andreianov, C. Cancès, and A. Moussa. A nonlinear time compactness result and applications to discretization of degenerate parabolic–elliptic PDEs. *J. Funct. Anal.*, 273(12):3633–3670, 2017.
- [3] J.-D. Benamou, G. Carlier, and M. Laborde. An augmented Lagrangian approach to Wasserstein gradient flows and applications. In *Gradient flows: from theory to application*, volume 54 of *ESAIM Proc. Surveys*, pages 1–17. EDP Sci., Les Ulis, 2016.
- [4] K. Brenner, C. Cancès, and D. Hilhorst. Finite volume approximation for an immiscible two-phase flow in porous media with discontinuous capillary pressure. *Comput. Geosci.*, 17(3):573–597, 2013.
- [5] C. Buet and S. Dellacherie. On the Chang and Cooper scheme applied to a linear Fokker-Planck equation. *Commun. Math. Sci.*, 8(4):1079–1090, 2010.
- [6] M. Burger, M. Di Francesco, J.-F. Pietschmann, and B. Schlake. Nonlinear cross-diffusion with size-exclusion. *SIAM J. Math. Anal.*, 46(6):2842–2871, 2010.
- [7] C. Cancès, C. Chainais-Hillairet, B. Merlet, F. Raimondi, and J. Venel. Mathematical analysis of a thermodynamically consistent reduced model for iron corrosion. working paper or preprint, January 2022.
- [8] C. Cancès, C. Chainais-Hillairet, J. Fuhrmann, and B. Gaudeul. A numerical analysis focused comparison of several finite volume schemes for a unipolar degenerated drift-diffusion model. *IMA J. Numer. Anal.*, 41(1):271–314, 2021.
- [9] C. Cancès, T. O. Gallouët, and G. Todeschi. A variational finite volume scheme for Wasserstein gradient flows. *Numer. Math.*, 146(3):437–480, 2020.
- [10] J. A. Carrillo, K. Craig, L. Wang, and C. Wei. Primal dual methods for Wasserstein gradient flows. *Found. Comput. Math.*, 2021. Online first.
- [11] M. Chatard. Asymptotic behavior of the Scharfetter-Gummel scheme for the drift-diffusion model. In *Finite volumes for complex applications VI. Problems & perspectives. Volume 1, 2*, volume 4 of *Springer Proc. Math.*, pages 235–243. Springer, Heidelberg, 2011.
- [12] K. Deimling. *Nonlinear functional analysis*. Springer-Verlag, Berlin, 1985.
- [13] R. Eymard, T. Gallouët, M. Ghilani, and R. Herbin. Error estimates for the approximate solutions of a nonlinear hyperbolic equation given by finite volume schemes. *IMA J. Numer. Anal.*, 18(4):563–594, 1998.
- [14] R. Eymard, T. Gallouët, C. Guichard, R. Herbin, and R. Masson. TP or not TP, that is the question. *Comput. Geosci.*, 18:285–296, 2014.
- [15] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. Ciarlet, P. G. (ed.) et al., in *Handbook of numerical analysis*. North-Holland, Amsterdam, pp. 713–1020, 2000.
- [16] K. Gärtner and L. Kamenski. Why Do We Need Voronoi Cells and Delaunay Meshes? In Vladimir A. Garanzha, Lennard Kamenski, and Hang Si, editors, *Numerical Geometry, Grid Generation and Scientific Computing*, Lecture Notes in Computational Science and Engineering, pages 45–60, Cham, 2019. Springer International Publishing.
- [17] M. Heida. Convergences of the squareroot approximation scheme to the Fokker–Planck operator. *Math. Models Methods Appl. Sci.*, 28(13):2599–2635, 2018.
- [18] M. Heida, M. Kantner, and A. Stephan. Consistency and convergence for a family of finite volume discretizations of the Fokker-Planck operator. *ESAIM: M2AN*, 55(6):3017–3042, 2021.
- [19] C. Kipnis, C. ; Landim. *Scaling limits of interacting particle systems*, volume 320. Springer, New-York, grundlehrender mathematischen Wissenschaften edition, 1999.
- [20] J. Leray and J. Schauder. Topologie et équations fonctionnelles. *Ann. Sci. École Norm. Sup.*, 51((3)):45–78, 1934.
- [21] W. Li, J. Lu, and L. Wang. Fisher information regularization schemes for Wasserstein gradient flows. *J. Comput. Phys.*, 416:109449, 2020.
- [22] H. C. Lie, K. Fackeldey, and M. Weber. A square root approximation of transition rates for a Markov state model. *SIAM J. Matrix Anal. Appl.*, 34(2):738–756, 2013.
- [23] D. Matthes and H. Osberger. Convergence of a variational Lagrangian scheme for a nonlinear drift diffusion equation. *ESAIM Math. Model. Numer. Anal.*, 48(3):697–726, 2014.
- [24] A. Mielke. A gradient structure for reaction-diffusion systems and for energy-drift-diffusion systems. *Nonlinearity*, 24(4):1329–1346, 2011.
- [25] A. Mielke, M. A. Peletier, and D. R. M. Renger. On the relation between gradient flows and the large-deviation principle, with applications to Markov chains and diffusion. *Potential Anal.*, 41(4):1293–1327, 2014.



- [26] M. A. Peletier, R. Rossi, G. Savaré, and O. Tse. Jump processes as generalized gradient flows. *Calc. Var. Partial Differential Equations*, 61:33, 2022.

CLÉMENT CANCÈS ([clement.cances@inria.fr](mailto:clement.cances@inria.fr)): UNIV. LILLE, INRIA, CNRS, UMR 8524 - LABORATOIRE PAUL PAINLEVÉ, F-59000 LILLE, FRANCE.

JULIETTE VENEL ([juliette.venel@uphf.fr](mailto:juliette.venel@uphf.fr)): UNIV. POLYTECHNIQUE HAUTS-DE-FRANCE, INSA HAUTS-DE-FRANCE, CERAMATHS – LABORATOIRE DE MATÉRIAUX CÉRAMIQUES ET DE MATHÉMATIQUES, F-59313 VALENCIENNES, FRANCE.