



**HAL**  
open science

# A convergent entropy diminishing finite volume scheme for a cross-diffusion system

Clément Cancès, Benoît Gaudeul

► **To cite this version:**

Clément Cancès, Benoît Gaudeul. A convergent entropy diminishing finite volume scheme for a cross-diffusion system. 2020. hal-02465431v1

**HAL Id: hal-02465431**

**<https://hal.science/hal-02465431v1>**

Preprint submitted on 3 Feb 2020 (v1), last revised 30 Jun 2020 (v2)

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

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

# A convergent entropy diminishing finite volume scheme for a cross-diffusion system

Clément Cancès\*

Benoît Gaudeul †

## Abstract

We study a two-point flux approximation finite volume scheme for a cross-diffusion system. The scheme is shown to preserve the key properties of the continuous systems, among which the decay of the entropy. The convergence of the scheme is established thanks to compactness properties based on the discrete entropy - entropy dissipation estimate. Numerical results illustrate the behavior of our scheme.

**Keywords.** Cross-diffusion system, Finite Volumes, discrete entropy method, convergence

**AMS subjects classification.** 35K51, 65M08, 65M12

## 1 Introduction

### 1.1 The system under study

The system studied in this paper has been originally introduced by [4] to model the production of solar panels using vapor deposition. In this system, we study the diffusion of  $N$  species whose respective concentrations are  $U = (u_1, \dots, u_N)$  in a (nonempty) connected bounded open domain  $\Omega$  of  $\mathbb{R}^d$  for a fixed time  $T$ . We denote by  $Q_T = (0, T) \times \Omega$ . The diffusion occurs through exchanges between different species which are quantified by the matrix  $A = (a_{i,j})$  of cross-diffusion coefficients. It leads to the following system of partial differential equations:

$$\partial_t u_i - \operatorname{div} \left( \sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) \right) = 0 \quad \text{in } Q_T \text{ for } i \in \llbracket 1, N \rrbracket. \quad (1)$$

The matrix  $A$  is assumed to be symmetric with nonnegative coefficients, i.e.  $a_{i,j} = a_{j,i} \geq 0$ .  $A$  does not depend on  $U$  and thus differs from the diffusion matrix  $D(U) = (d_{i,j}(U))$  defined by

$$d_{i,j}(U) = \delta_{i,j} \sum_{k \neq i} a_{i,k} u_k - a_{i,j} u_i,$$

where  $\delta_{i,j}$  stands for Kronecker symbol, such that the problem (1) rewrites

$$\partial_t U - \operatorname{div}(D(U)\nabla U) = 0. \quad (2)$$

---

\*Inria, Univ. Lille, CNRS, UMR 8524 - Laboratoire Paul Painlevé, F-59000 Lille (clement.cances@inria.fr)

†Univ. Lille, CNRS, UMR 8524, Inria - Laboratoire Paul Painlevé, F-59000 Lille (benoit.gaudeul@univ-lille.fr)

System (2) enters the family of the nonlinear cross-diffusion systems since  $D$  depends on  $U$  and has nonzero off-diagonal entries. Challenges both from the analytical and numerical points of view come from the presence of off-diagonal zeros in  $A$ . In the previous contributions [6, 23, 7], the zeros are integrated through the assumption that the cross-diffusion occurs with and only with a solvent specie. Until Section 4 we will not make any assumption about the zeros of  $A$ . A non-degeneracy assumption will be further assumed on Section 4, but our convergence result could extend to the particular cross-diffusion matrices considered in [23, 7, 22].

We supplement system (1) with no-flux boundary conditions

$$\sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) \cdot n = 0 \quad \text{on } (0, T) \times \partial\Omega, \quad i \in \llbracket 1, N \rrbracket. \quad (3)$$

The initial concentration  $U^0 = (u_1^0, \dots, u_N^0)$  is supposed to be measurable and to map  $\Omega$  into

$$\mathcal{A} = \left\{ U = (u_1, \dots, u_N) \in \mathbb{R}_+^N \left| \sum_{i=1}^N u_i = 1 \right. \right\},$$

so we write in the condensed form  $U^0 \in L^\infty(\Omega, \mathcal{A})$ . Finally, we assume that all the chemical species under consideration are present:

$$\int_{\Omega} u_i^0 dx > 0, \quad \forall i \in \llbracket 1, N \rrbracket. \quad (4)$$

## 1.2 Formal structure

This system has several structural properties, the goal of this subsection is to exhibit them. The calculations presented in this section are formal: we assume that the solutions to (1) enjoy enough regularity to justify the calculations below. Rigorous proofs at the continuous level for the system under consideration here can be found in [4, 5] (see also [23]). The properties listed here can also be obtained by passing to the limit in the numerical scheme. The first property we point out is the conservation of mass for all the species involved in System (1).

**Lemma 1.1 (conservation of mass)** (1) and (3) corresponding to an initial data  $U^0 \in L^\infty(\Omega, \mathcal{A})$ , then

$$\int_{\Omega} u_i(t, x) dx = \int_{\Omega} u_i^0(x) dx, \quad \forall t \in [0, T], \forall i \in \llbracket 1, N \rrbracket.$$

*Proof.* Let  $U$  be a solution of (1),  $t \in [0, T]$ ,  $i \in \llbracket 1, N \rrbracket$ , and let  $\varphi(x, s) = \mathbf{1}_{[0, t]}(s)$ . With this particular choice of  $\varphi$ , we have for all  $s$  that

$$\int_{\Omega} \operatorname{div} \left( \sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) \right) \varphi(x, s) dx = - \int_{\Omega} \sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) \nabla \varphi(x, s) dx = 0.$$

Hence, using  $\varphi$  as a test function in (1), we have:

$$\int_0^t \frac{d}{ds} \left( \int_{\Omega} u_i(x, s) dx \right) ds = 0.$$

The fundamental theorem of calculus yields the desired lemma.  $\square$

The symmetry of the matrix  $A = (a_{i,j})$  yields:

$$\sum_{i=1}^N \sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) = 0.$$

Therefore, a solution  $U$  to (1) satisfies  $\partial_t \sum_{i=1}^N u_i = 0$ . Admit that  $u_i(t,x) \geq 0$  for all  $t > 0$  (this will be proved in the discrete setting and is proved in [5, Proposition 2.2] in the continuous setting), then the admissibility condition encoded in  $\mathcal{A}$  is preserved along time.

**Lemma 1.2** *Let  $U$  be a solution to (1) and (3) corresponding to an initial data  $U^0 \in L^\infty(\Omega; \mathcal{A})$ , then  $U(t,x) \in \mathcal{A}$  for all  $(t,x) \in \mathcal{A}$ , i.e.,  $U \in L^\infty(Q_T; \mathcal{A})$ .*

The next property we want to highlight at the continuous level is the decay of entropy. Using the chain rule  $\nabla c = c \nabla \ln(c)$ , the system (1) rewrites

$$\partial_t u_i - \operatorname{div} \left( \sum_{j=1}^N a_{i,j} u_i u_j (\nabla \ln(u_i) - \nabla \ln(u_j)) \right) = 0, \quad i \in \llbracket 1, N \rrbracket. \quad (5)$$

**Proposition 1.3** *Introduce the functional:*

$$E : U \mapsto \int_{\Omega} \sum_{i=1}^N u_i \ln(u_i) dx$$

mapping  $L^\infty(\Omega; \mathcal{A})$  into  $\mathbb{R}$ , then  $E$  is a Lyapunov functional for the system (3)–(5). More precisely, the following entropy - entropy dissipation estimate holds:

$$\frac{d}{dt} E(U) + \int_{\Omega} \left( \sum_{1 \leq i < j \leq N} a_{i,j} u_i u_j |\nabla \ln(u_i) - \nabla \ln(u_j)|^2 \right) dx = 0. \quad (6)$$

*Proof.* First, we notice that thanks to the conservation of mass:

$$\frac{d}{dt} E(U) = \frac{d}{dt} \int_{\Omega} \sum_{i=1}^N u_i (\ln(u_i) - 1) = \int_{\Omega} \sum_{i=1}^N \ln(u_i) \partial_t u_i.$$

Then multiply Equation (5) by  $\ln(u_i)$  and integrate by part in order to get:

$$\int_{\Omega} \ln(u_i) \partial_t u_i + \int_{\Omega} \left( \sum_{j=1}^N a_{i,j} u_i u_j \nabla \ln(u_i) \cdot (\nabla \ln(u_i) - \nabla \ln(u_j)) \right) = 0.$$

Summing over  $i \in \llbracket 1, N \rrbracket$  yields the announced result thanks to the symmetry of  $A$ .  $\square$

The entropy - entropy dissipation relation (6) is key in the analysis of many cross-diffusion systems, as exposed in [25, 26]. It will also play a central role in this paper. Assume that

$$\min_{i \neq j} a_{i,j} > 0 \quad (7)$$

as it will be done in Section 4. As a consequence of the inequality

$$\sum_{i=1}^N \int_{\Omega} |\nabla u_i|^2 \leq 4 \sum_{i=1}^N \int_{\Omega} |\nabla \sqrt{u_i}|^2 \leq \frac{1}{\min_{i \neq j} a_{i,j}} \int_{\Omega} \sum_{1 \leq i < j \leq N} a_{i,j} u_i u_j |\nabla \ln(u_i) - \nabla \ln(u_j)|^2,$$

we deduce from (6) a  $L^2(0, T; H^1(\Omega))^N$  estimate on  $U$ . This motivates the following notion of weak solution.

**Definition 1.4** A weak solution  $U$  to (1) and (3) corresponding to the initial profile  $U^0 \in L^\infty(\Omega; \mathcal{A})$  is a function of  $L^\infty(Q_T; \mathcal{A}) \cap L^2([0, T]; H^1(\Omega))^N$  satisfying,  $\forall i \in \llbracket 1, N \rrbracket$ ,  $\forall \varphi \in C_c^\infty([0, T] \times \overline{\Omega})$ :

$$\iint_{Q_T} u_i \partial_t \varphi dx dt - \int_{\Omega} u_i^0 \varphi(0, \cdot) dx + \iint_{Q_T} \sum_{j=1}^N a_{i,j} (u_j \nabla u_i - u_i \nabla u_j) \nabla \varphi = 0. \quad (8)$$

The regularity requirement on a weak solution  $U$  is natural in the setting where Assumption (7) holds. In this case, the solution even enjoys a stronger regularity typically, as established in the recent contribution [5]. In the case where (7) is not fulfilled (but under a structural assumption on the matrix  $A$ ), a more involved notion of weak solution has to be introduced, cf. [23].

There is an important property that relates the model (1) to classical Fickian diffusion. As a consequence of Lemma 1.2, one can rewrite

$$\operatorname{div} \left( \sum_{j=1}^N (u_j \nabla u_i - u_i \nabla u_j) \right) = \Delta u_i, \quad i \in \llbracket 1, N \rrbracket. \quad (9)$$

As a consequence, if all the  $a_{i,j}$  are equal to some  $a \in \mathbb{R}$ , then the system (1) reduces to  $N$  uncoupled heat equations  $\partial_t u_i = a \Delta u_i$ . Based on the identity (9), we can rewrite the system (1) under the form

$$\partial_t u_i - a^* \Delta u_i - \operatorname{div} \left( \sum_{j=1}^N (a_{i,j} - a^*) (u_j \nabla u_i - u_i \nabla u_j) \right) = 0, \quad i \in \llbracket 1, N \rrbracket, \quad (10)$$

where  $a^* \in \mathbb{R}$  is arbitrary for the moment. The formulation (10) is at the basis of our discretization.

### 1.3 Objectives

The goal of this paper is to build and analyze a numerical scheme preserving the properties discussed in the previous section, namely:

- the non-negativity of the concentrations;
- the conservation of mass (Lemma 1.1);
- the preservation of the volume filling constraint (Lemma 1.2);
- the entropy-entropy dissipation relation (Proposition 1.3).

The construction of our scheme is the purpose of Section 2. In Section 3, we will show the existence of solutions to this scheme and the preservation of discrete counterparts to the previously listed physical properties. Section 4 is devoted to the convergence of the numerical scheme toward weak solutions provided Assumption (7) is satisfied. Finally in Section 5, we show the outcomes of some numerical experiments.

Before entering the core of the paper, let us mention that the development of numerical analysis for cross-diffusion systems is quite recent. To our knowledge, the first convergence study of a finite volume approximation for a non-degenerate cross-diffusion problem was carried out in [3]. This contributions is based on classical quadratic energy estimate. The implementation of the discrete entropy method [9] for cross-diffusion systems is more recent. Let us cite [1, 2] where upstream mobility finite volume and control volume finite elements schemes for a multiphase extension of the porous medium equation are

studied. Upwinding is also used in [7] to approximate the solution of a system which is very close to the problem (1) under study, or in [8] for a problem in which nonlocal interactions are also considered. As a consequence of the upwind choice for the mobility, the schemes presented in [1, 2, 7] and [8] are first order accurate in space. An natural solution to pass to order two is to rather consider mobilities given by arithmetic means [12]. The motivation of the finite element scheme proposed in [27] is also the same. However, the scheme proposed in [27] is expressed in entropy (or dual) variables (in our context  $\log(u_i)$ ) leading to computational difficulties when the concentrations are close to 0. Finally, let us mention the extension to higher order discontinuous Galerkin methods proposed in [29].

## 2 Finite Volume approximation

This section is organized as follows. First, in Section 2.1, we state the requirements on the mesh and fix some notations. Then in Section 2.2, we describe the numerical scheme to be studied in this paper. It is based on Formulation (10) of the problem. Then in Section 2.3, we state our two main results. The first one, namely Theorem 2.2, focuses on the case of a fixed mesh. We are interested in the existence of a solution to the nonlinear system corresponding to the scheme, and the dissipation of the entropy at the discrete level. More precisely, one establishes that the studied scheme satisfies a discrete entropy - entropy dissipation inequality that can should be thought as a counterpart to Proposition 1.3. Our second main result, namely Theorem 2.3, is devoted to the convergence of the scheme towards a weak solution as the time step and the mesh size tend to 0.

### 2.1 Discretization of $(0, T) \times \Omega$

The scheme we propose relies on two-point flux approximation (TPFA) finite volumes. As explained in [14, 18, 21], this approach appears to be very efficient as soon as the continuous problem to be solved numerically is isotropic and one has the freedom to choose a suitable mesh fulfilling the so-called orthogonality condition [24, 19]. We recall here the definition of such a mesh, which is illustrated in Figure 1.

**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*

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

- (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  $\bar{K} \cap \bar{L}$  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  $(x_K)_{K \in \mathcal{T}}$  satisfy  $x_K \in K$ , and are such that, if  $K, L \in \mathcal{T}$  share a face  $K|L$ , then the vector  $x_L - x_K$  is orthogonal to  $K|L$ .*

- (iv) For the boundary faces  $\sigma \subset \partial\Omega$ , we assume that either  $\sigma \subset \Gamma_D$  or  $\sigma \subset \bar{\Gamma}_N$ . For  $\sigma \subset \partial\Omega$  with  $\sigma \in \mathcal{E}_K$  for some  $K \in \mathcal{T}$ , we assume additionally that there exists  $x_\sigma \in \sigma$  such that  $x_\sigma - x_K$  is orthogonal to  $\sigma$ .

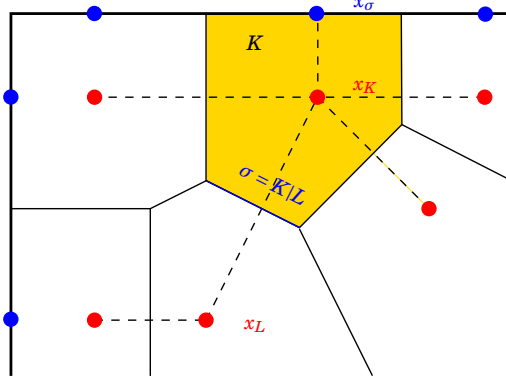


Figure 1: Illustration of an admissible mesh as in Definition 2.1.

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  $\mathcal{E}_{K,\text{int}} = \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$  the set of its faces that belong to  $\mathcal{E}_{\text{int}}$ . For such a face  $\sigma \in \mathcal{E}_{K,\text{int}}$ , we may write  $\sigma = K|L$ , meaning that  $\sigma = \bar{K} \cap \bar{L}$ , where  $L \in \mathcal{T}$ .

Given  $\sigma \in \mathcal{E}$ , we let

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

We finally introduce the size  $h_{\mathcal{T}}$  and the regularity  $\zeta_{\mathcal{T}}$  (which is assumed to be positive) of a discretization  $(\mathcal{T}, \mathcal{E}, (x_K)_{K \in \mathcal{T}})$  of  $\Omega$  by setting

$$h_{\mathcal{T}} = \max_{K \in \mathcal{T}} \text{diam}(K), \quad \zeta_{\mathcal{T}} = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d(x_K, \sigma)}{d_\sigma}.$$

Concerning the time discretization of  $(0, T)$ , we consider an increasing finite family of times  $0 = t_0 < t_1 < \dots < t_{N_T} = T$ . We denote by  $\Delta t_n = t_n - t_{n-1}$  for  $n \in \{1, \dots, N_T\}$ , by  $\mathbf{\Delta t} = (\Delta t_n)_{1 \leq n \leq N_T}$ , and by  $h_T = \max_{1 \leq n \leq N_T} \Delta t_n$ . In what follows, we will use boldface notations for mesh-indexed families, typically for elements of  $\mathbb{R}^{\mathcal{T}}$ ,  $(\mathbb{R}^{\mathcal{T}})^N$ ,  $(\mathbb{R}^{\mathcal{T}})^{N_T}$ , or even  $(\mathbb{R}^{\mathcal{T}})^{N \times N_T}$ .

## 2.2 Numerical scheme

The initial data  $U^0 \in L^\infty(\Omega; \mathcal{A})$  is discretized into

$$\mathbf{U}^0 = (\mathbf{u}_i^0)_{i \in \llbracket 1, N \rrbracket} \in (\mathbb{R}^{\mathcal{T}})^N = \left( u_{i,K}^0 \right)_{K \in \mathcal{T}, i \in \llbracket 1, N \rrbracket}$$

by setting

$$u_{i,K}^0 = \frac{1}{m_K} \int_K u_i^0(x) dx, \quad \forall K \in \mathcal{T}, i \in \llbracket 1, N \rrbracket. \quad (11)$$

Assume that  $\mathbf{U}^{n-1} = (u_{i,K}^{n-1})_{K \in \mathcal{T}, i \in \llbracket 1, N \rrbracket}$  is given for some  $n \geq 1$ , then we have to define how to compute  $\mathbf{U}^n = (u_{i,K}^n)_{K \in \mathcal{T}, i \in \llbracket 1, N \rrbracket}$ .

First, we introduce some notations. Given any discrete scalar field  $\mathbf{c} = (c_K)_{K \in \mathcal{T}} \in \mathbb{R}^{\mathcal{T}}$ , we define for all cell  $K \in \mathcal{T}$  and interface  $\sigma \in \mathcal{E}_K$  the mirror value  $c_{K\sigma}$  of  $c_K$  across  $\sigma$  by setting:

$$c_{K\sigma} = \begin{cases} c_L & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ c_K & \text{if } \sigma \in \mathcal{E}_{\text{ext}}. \end{cases} \quad (12)$$

We also define the oriented and absolute jumps of  $\mathbf{c}$  across any edge by

$$D_{K\sigma} \mathbf{c} = c_{K\sigma} - c_K, \quad D_\sigma \mathbf{c} = |D_{K\sigma} \mathbf{c}|, \quad \forall K \in \mathcal{T}, \forall \sigma \in \mathcal{E}_K.$$

The scheme is based on the formulation (10). It requires the introduction of a parameter  $a^\star$  on which we only have the following requirements:

$$a^\star > 0 \quad \text{and} \quad a^\star \geq \min_{(i,j)} a_{i,j}. \quad (13)$$

The conservation laws are discretized in a conservative way with a time discretization relying on the backward Euler scheme:

$$m_K \frac{u_{i,K}^n - u_{i,K}^{n-1}}{\Delta t_n} + \sum_{\sigma \in \mathcal{E}_K} F_{i,K\sigma}^n = 0, \quad \forall K \in \mathcal{T}, \forall i \in \llbracket 1, N \rrbracket. \quad (14a)$$

The discrete fluxes are computed thanks to a formula based on (10) and on TPFA finite volumes:

$$F_{i,K\sigma}^n = -a^\star \tau_\sigma D_{K\sigma} \mathbf{u}_i^n - \tau_\sigma \left( \sum_{j=1}^N (a_{i,j} - a^\star) (u_{j,\sigma}^n D_{K\sigma} u_i^n - u_{i,\sigma}^n D_{K\sigma} u_j^n) \right), \quad (14b)$$

for all  $K \in \mathcal{T}$ ,  $\sigma \in \mathcal{E}_K$  and  $i \in \llbracket 1, N \rrbracket$ . Edge values  $(u_{j,\sigma}^n)_j$  of the concentrations  $u_j$  appears in Formula (14b). It is deduced from  $u_{j,K}^n$  and  $u_{j,K\sigma}^n$  thanks to a logarithmic mean, i.e.,

$$u_{j,\sigma}^n = \begin{cases} 0 & \text{if } \min(u_{j,K}^n, u_{j,K\sigma}^n) \leq 0, \\ u_{j,K}^n & \text{if } 0 \leq u_{j,K}^n = u_{j,K\sigma}^n, \\ \frac{u_{j,K}^n - u_{j,K\sigma}^n}{\ln(u_{j,K}^n) - \ln(u_{j,K\sigma}^n)} & \text{otherwise.} \end{cases} \quad (14c)$$

This choice for the edge concentration is crucial for the preservation at the discrete level of a discrete entropy - entropy dissipation inequality similar to the one highlighted in Proposition 1.3. Equations (14b) and (12) implies that for all  $\sigma \in \mathcal{E}_{\text{ext}}$ :  $F_{i,K\sigma}^n = 0$ , so that the no-flux boundary condition (3) is taken into account.



**Remark 2.1** Let us highlight why the choice of a strictly positive  $a^*$  is important. Consider a mesh with two cells  $K, L$ , and one edge. We consider two species and let  $u_K^0 = (0, 1)$  and  $u_L^0 = (1, 0)$ . We have:  $u_{1,K|L}^0 = 0$  and  $u_{2,K|L}^0 = 0$ , hence, if  $a^* = 0$ , the initial condition is a stationary solution even though this is not expected for a discretization of the heat equation. Setting  $a^* > 0$  eliminates these spurious solutions. The choice of  $a^*$  has a strong influence on the numerical outcomes, as it will be shown in Section 5, but we don't have a clear understanding yet on the methodology to choose an optimal  $a^*$ . What is clear is that  $a^*$  has to be chosen in the interval  $[\min_{i \neq j} a_{i,j}, \max_{i \neq j} a_{i,j}]$ . A tentative non-optimal formula is proposed in Section 5.

### 2.3 Main results and organization

The first theorem proven in this paper concerns the existence of discrete solutions for a given mesh, and the preservation of the structural properties listed in Section 1.3:

- the mass of each specie is conserved along the time steps;
- the concentrations are (strictly) positive and sum to 1 in all the cells, i.e.,  $U_K^n \in \mathcal{A}$  for all  $K \in \mathcal{T}$  and  $n \geq 1$ ;
- the discrete counterpart of the entropy decays along time.

For this last property, we need to introduce the discrete entropy functional  $E_{\mathcal{T}}$ , which is defined by:

$$E_{\mathcal{T}}(\mathbf{U}) = \sum_{K \in \mathcal{T}} \sum_{i=1}^N m_K u_{i,K} \ln u_{i,K}, \quad \forall \mathbf{U} = (u_{i,K})_{K \in \mathcal{T}, i \in \llbracket 1, N \rrbracket} \in \mathcal{A}^{\mathcal{T}}. \quad (15)$$

As stated in Theorem 2.2 below, the nonlinear system corresponding to our scheme (14) admits solutions which preserve the physical bounds on the concentrations and the decay of the entropy.

**Theorem 2.2** Let  $(\mathcal{T}, \mathcal{E}, (x_K)_{K \in \mathcal{T}})$  be an admissible mesh and let  $\mathbf{U}^0$  be defined by (11). Then, for all  $1 \leq n \leq N_T$ , the nonlinear system of equations (12) – (14), has a positive solution  $\mathbf{U}^n \in \mathcal{A}^{\mathcal{T}}$ . Moreover, such a solution satisfies  $E_{\mathcal{T}}(\mathbf{U}^n) \leq E_{\mathcal{T}}(\mathbf{U}^{n-1})$  for all  $n \in \llbracket 1, N_T \rrbracket$ ,  $\sum_{K \in \mathcal{T}} m_K u_{i,K}^n = \int_{\Omega} u_i^0$  for all  $i \in \llbracket 1, N \rrbracket$  and  $n \in \llbracket 0, N_T \rrbracket$ .

The proof of Theorem 2.2 will be the purpose of Section 3. With a discrete solution  $(\mathbf{U}^n)_{1 \leq n \leq N_T}$  to the scheme (14) at hand, we can define the piecewise constant approximate solution  $U_{\mathcal{T}, \Delta t} = (u_{i, \mathcal{T}, \Delta t})_{i \in \llbracket 1, N \rrbracket} : Q_T \rightarrow \mathcal{A}$  defined almost everywhere by

$$U_{\mathcal{T}, \Delta t}(t, x) = U_K^n \quad \text{if } (t, x) \in (t_{n-1}, t_n] \times K.$$

This definition will be developed in Section 4 and supplemented by other reconstruction operators. Let  $(\mathcal{T}_m, \mathcal{E}_m, (x_K)_{K \in \mathcal{T}_m})_{m \geq 1}$  be a sequence of admissible discretizations with  $h_{\mathcal{T}_m}, h_{T,m}$  tending to 0 as  $m$  tends to  $+\infty$ , while the regularity  $\zeta_{\mathcal{T}_m}$  remains uniformly bounded from below by a positive constant  $\zeta^*$ . Thanks to Theorem 2.2, we dispose of a family  $\mathbf{U}_m$  of solutions to our scheme. The convergence of  $\mathbf{U}_m$  is the purpose of Theorem 2.3 whose proof is detailed in Section 4.

**Theorem 2.3** Assume that the nondegeneracy assumption (7) holds. Given any sequence of solutions  $\mathbf{U}_m = (u_{i,K}^n)_{i \in \llbracket 1, N \rrbracket, K \in \mathcal{T}_m, 1 \leq n \leq N_{T,m}}$ , there exists at least one  $U \in L^\infty(Q_T; \mathcal{A}) \cap L^2((0, T); H^1(\Omega))$  such that, up to a subsequence,

$$U_{\mathcal{T}_m, \Delta t_m} \xrightarrow{m \rightarrow \infty} U \quad \text{strongly in } L^p(Q_T), \text{ for any } 1 \leq p < \infty, \quad (16)$$

Moreover,  $U$  is a weak solution in the sense of Definition 1.4.

### 3 Numerical analysis on a fixed mesh

This section is devoted to the proof of Theorem 2.2. In Section 3.1, we establish a priori estimates on an slightly modified scheme that will be shown to reduce to the original scheme (14). Then in Section 3.2, we apply a topological degree argument to prove the existence of solutions to our scheme. Section 3.3 is devoted to the proof of the entropy dissipation property.

To prove the existence of solutions to the system of equations (14), we need the inequality  $\sum_i u_{i,\sigma} \leq 1$ . We then slightly modify (14) by adding the following equation:

$$\widetilde{u}_{i,\sigma}^n = \frac{u_{i,\sigma}^n}{\max(1, \sum_{j=1}^N u_{j,\sigma}^n)},$$

and replacing  $u_{i,\sigma}^n$  by  $\widetilde{u}_{i,\sigma}^n$  in (14b). We will denote this new system (S) and see in Proposition 3.4 that its solutions satisfy  $\sum_i u_{i,\sigma} \leq 1$ , so that  $\widetilde{u}_{i,\sigma}^n = u_{i,\sigma}^n$ . Whence they also satisfy the original system of equations.

#### 3.1 A priori estimates

The first lemma shows the nonnegativity of the solutions to (S).

**Lemma 3.1** *Given a nonnegative  $U^{n-1}$ , any solution  $U^n$  to (S) is also nonnegative.*

*Proof.* Let  $U^n$  be a solution of (S) and let  $i \in \llbracket 1, N \rrbracket$ . We consider a cell  $K \in \mathcal{T}$  where  $u_i^n$  reaches its minimum, i.e.,  $u_{i,K}^n \leq u_{i,L}^n$  for all  $L \in \mathcal{T}$ , and assume for contradiction that  $u_{i,K}^n$  is (strictly) negative. Equation (14b) then gives:

$$m_K \frac{u_{i,K}^n - u_{i,K}^{n-1}}{\Delta t_n} = - \sum_{\sigma \in \mathcal{E}_K} F_{K\sigma}^n.$$

The term on the left hand side is negative since  $u_{i,K}^{n-1} \geq 0 > u_{i,K}^n$ , whereas the right-hand side may be simplified noticing that  $\widetilde{u}_{i,\sigma}^n = 0$ :

$$\sum_{\sigma \in \mathcal{E}_K} a^* \tau_\sigma D_{K\sigma} u_i^n + \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \sum_{j=1}^N (a_{i,j} - a^*) \widetilde{u}_{j,\sigma}^n D_{K\sigma} u_i^n = - \sum_{\sigma \in \mathcal{E}_K} F_{K\sigma}^n < 0.$$

Noticing that  $D_{K\sigma} u_i^n \geq 0$ ,  $\widetilde{u}_{j,\sigma}^n \geq 0$ , and  $\sum_{j=1}^N \widetilde{u}_{j,\sigma}^n \leq 1$  we obtain that

$$0 \leq \sum_{\sigma \in \mathcal{E}_K} a^* (1 - \sum_{j=1}^N \widetilde{u}_{j,\sigma}^n) \tau_\sigma D_{K\sigma} u_i^n < 0,$$

which is absurd, hence the desired result. □

Let us now show that the concentrations sum to 1 in all the cells.

**Lemma 3.2** *Given  $U^{n-1}$  in  $\mathcal{A}^{\mathcal{T}}$ , any solution  $U^n$  to (S) is also in  $\mathcal{A}^{\mathcal{T}}$ .*

*Proof.* Thanks to Lemma 3.1, it suffices to show that  $\sum_{i=1}^N u_{i,K}^n = 1$  for all  $K \in \mathcal{F}$ . Let  $\mathbf{U}^n$  be a solution to (S). Using (14b) in (14a) and summing over the species leads to:

$$\begin{aligned} \frac{\sum_{i=1}^N u_{i,K}^n - \sum_{i=1}^N u_{i,K}^{n-1}}{\Delta t_n} m_K - a^* \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma D_{K\sigma} \sum_i \mathbf{u}_i \\ - \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \sum_i \left( \sum_{j=1}^N (a_{i,j} - a^*) \left( \widetilde{u}_{j,\sigma}^n D_{K\sigma} \mathbf{u}_i - \widetilde{u}_{i,\sigma}^n D_{K\sigma} \mathbf{u}_j \right) \right) = 0, \quad \forall K \in \mathcal{F}. \end{aligned}$$

The third term of the left-hand side vanishes thanks to the symmetry of  $A$ , so that

$$\frac{\sum_{i=1}^N u_{i,K}^n - \sum_{i=1}^N u_{i,K}^{n-1}}{\Delta t_n} m_K - a^* \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma D_{K\sigma} \sum_i \mathbf{u}_i = 0, \quad \forall K \in \mathcal{F}.$$

The discrete quantity  $\sum_i \mathbf{u}_i$  is solution to the classical backward Euler TPFA scheme for the heat equation, which is well posed. So  $\sum_i \mathbf{u}_i^n = \sum_i \mathbf{u}_i^{n-1} = \mathbf{1}$  is its unique solution, hence the desired result.  $\square$

## 3.2 Existence of solutions

Using the tools exposed in the previous subsection, we may derive the existence of a solution to (S):

**Proposition 3.3** *Given  $\mathbf{U}^{n-1}$  in  $\mathcal{A}^{\mathcal{F}}$ , there exists at least one solution to (S) in  $\mathcal{A}^{\mathcal{F}}$ .*

*Proof.* The proof relies on a topological degree argument [28, 13]. The idea is to transform continuously our complex nonlinear system into a linear system while guaranteeing that the a priori estimates controlling the solution remain valid all along the homotopy. We sketch the main ideas of the proof, making the homotopy explicit. We are interested in the existence of zeros for a functional

$$\mathcal{H}: \begin{cases} [0, 1] \times (\mathbb{R}^N)^{\mathcal{F}} \rightarrow (\mathbb{R}^N)^{\mathcal{F}} \\ (\lambda, \mathbf{U}) \mapsto \mathcal{H}(\lambda, \mathbf{U}) \end{cases}$$

that boils down to the scheme (S) when  $\lambda = 1$ . In our case, we set:

$$\begin{aligned} \mathcal{H}(\lambda, \mathbf{U})_{i,K} = \frac{u_{i,K} - u_{i,K}^{n-1}}{\Delta t_n} m_K - a^* \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma D_{K\sigma} \mathbf{u}_i \\ - \lambda \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \left( \sum_{j=1}^N (a_{i,j} - a^*) \left( \widetilde{u}_{j,\sigma} D_{K\sigma} \mathbf{u}_i - \widetilde{u}_{i,\sigma} D_{K\sigma} \mathbf{u}_j \right) \right), \quad \forall K \in \mathcal{F}, \forall i \in \llbracket 1, N \rrbracket. \end{aligned} \quad (17)$$

One notices that  $\mathcal{H}(0, \mathbf{U}) = \mathbf{0}$  is the classical heat equation, the solution of which belongs to  $\mathcal{A}^{\mathcal{F}}$ . Therefore, fixing  $\eta > 0$ , the relatively compact open set

$$\mathcal{A}_\eta^{\mathcal{F}} = \left\{ \mathbf{U} \in \mathbb{R}^{\mathcal{F}} \mid \inf_{\mathbf{V} \in \mathcal{A}^{\mathcal{F}}} \|\mathbf{U} - \mathbf{V}\| < \eta \right\}$$

has a topological degree equal to 1. Note that the choice of the norm in the definition of  $\mathcal{A}_\eta^{\mathcal{F}}$  is not important since the dimension is finite. Moreover, thanks to Lemma 3.2, the solutions  $\mathbf{u}^{(\lambda)}$  of

$\mathcal{H}(\lambda, \mathbf{U}) = \mathbf{0}$  remains in  $\mathcal{A}^{\mathcal{F}}$ , thus in the interior of  $\mathcal{A}_\eta^{\mathcal{F}}$ . Thus the topological degree of  $\mathcal{A}_\eta^{\mathcal{F}}$  for  $\lambda = 1$  is still equal to 1, hence the existence of (at least) one solutions to (S). Since  $\eta > 0$  is arbitrary, then there is a solution in  $\mathcal{A}^{\mathcal{F}} = \bigcap_{\eta > 0} \mathcal{A}_\eta^{\mathcal{F}}$ .  $\square$

To prove the Theorem 2.2, we need to transfer this result on the original system.

**Proposition 3.4** *A solution  $\mathbf{U}^n$  of (S) is a solution of (14). Reciprocally, a solution of (14) in  $\mathcal{A}^{\mathcal{F}}$  is a solution of (S).*

*Proof.* Let  $\mathbf{U}^n$  be a solution of (S). A simple convexity argument shows that the logarithmic mean of two nonnegative number is smaller than the arithmetic mean, so that  $u_{i,\sigma}^n \leq \frac{u_{i,K}^n + u_{i,K\sigma}^n}{2}$ . Summing w.r.t.  $i \in \llbracket 1, N \rrbracket$  and using that the solution  $\mathbf{U}$  of (S) belongs to  $\mathcal{A}^{\mathcal{F}}$ , one gets that  $\sum_i u_{i,\sigma}^n \leq 1$  for all  $\sigma \in \mathcal{E}$ . Therefore  $\widetilde{u}_{i,\sigma}^n = u_{i,\sigma}^n$  and  $\mathbf{U}^n$  is a also solution to (14). The proof of the reverse implication follows the same lines.  $\square$

### 3.3 Entropy dissipation

We intend here to prove a discrete counterpart to Proposition 1.3. The proof will be very similar and requires a discrete counterpart of the conservation of mass (Lemma 1.1).

**Lemma 3.5** *Given any  $\mathbf{U}^{n-1} \in \mathcal{A}^{\mathcal{F}}$ , any solution  $\mathbf{U}^n$  to (14) satisfies:*

$$\sum_{K \in \mathcal{T}} m_K u_{i,K}^n = \sum_{K \in \mathcal{T}} m_K u_{i,K}^{n-1} = \int_{\Omega} u_i^0 dx, \quad \forall i \in \llbracket 1, N \rrbracket.$$

The proof of this lemma is a straightforward calculation based on equation (14a), the conservativity of the fluxes, and the definition (11) of the discrete initial condition. With this lemma and Proposition 3.4, we can refine the result Lemma 3.1 to get the strict positivity of any solution to (14) belonging to  $\mathcal{A}^{\mathcal{F}}$ .

**Lemma 3.6** *Let  $\mathbf{U}^{n-1} \in \mathcal{A}^{\mathcal{F}}$  be such that  $\sum_K m_K u_{i,K}^{n-1} > 0$  for all  $i \in \llbracket 1, N \rrbracket$ , then any solution to (14) in  $\mathcal{A}^{\mathcal{F}}$  is positive:  $u_{i,K}^n > 0$  for all  $i \in \llbracket 1, N \rrbracket$  and all  $K \in \mathcal{T}$ .*

*Proof.* Let  $\mathbf{U}^n \in \mathcal{A}^{\mathcal{F}}$  be a solution to the scheme (14), and let  $i \in \llbracket 1, N \rrbracket$ . We know from Lemma 3.1 that  $\mathbf{u}_i^n \geq \mathbf{0}$ . Assume for contradiction that there exists one cell  $K$  such that  $u_{i,K}^n$  vanishes. Using Lemma 3.5 and the connectivity of  $\Omega$ , there exists  $\sigma = K|L \in \mathcal{E}^{\text{int}}$  such that  $u_{i,K}^n = 0$  and  $u_{i,L}^n > 0$ . Then  $u_{i,\sigma}^n = 0$  and as in the proof of Lemma 3.1:

$$a^* \left(1 - \sum_{j=1}^N u_{j,\sigma}^n\right) \tau_\sigma D_{K\sigma} \mathbf{u}_i^n \leq 0.$$

Using  $u_{j,\sigma}^n \leq \frac{u_{j,K}^n + u_{j,L}^n}{2}$  and  $u_{i,\sigma}^n = 0$  we deduce that

$$\sum_{j=1}^N u_{j,\sigma}^n \leq \sum_{j \neq i} \frac{u_{j,K}^n + u_{j,L}^n}{2} \leq 1 - \frac{u_{i,L}^n}{2} < 1.$$

Therefore  $a^* \left(1 - \sum_{j=1}^N \widetilde{u}_{j,\sigma}^n\right) \tau_\sigma > 0$ , and since  $D_{K\sigma} \mathbf{u}_i^n > 0$ , we deduce that:

$$0 < a^* \left(1 - \sum_{j=1}^N \widetilde{u}_{j,\sigma}^n\right) \tau_\sigma D_{K\sigma} \mathbf{u}_i^n \leq 0.$$

As this statement is absurd, our assumption was false, hence the desired result.  $\square$  As in the continuous case, we will use the conservation of mass (Lemma 3.5) and a discrete equivalent of the chain rule  $\nabla c = c \nabla \ln c$ . This equivalent writes

$$D_{K\sigma} \mathbf{u}_i^n = u_{i,\sigma}^n D_{K\sigma} \ln(\mathbf{u}_i^n), \quad \forall i \in \llbracket 1, N \rrbracket, \forall K \in \mathcal{F}. \quad (18)$$

The above discrete chain rule follows from the definition (14c) of  $u_{i,\sigma}^n$  and the positivity of solutions to (14) which gives a sense to  $\ln(\mathbf{u}_i^n)$ .

Using (18) in (14b),  $\mathbf{U}^n$  satisfies

$$\begin{aligned} \frac{u_{i,K}^n - u_{i,K}^{n-1}}{\Delta t_n} m_K - \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \left( \sum_{j=1}^N (a_{i,j} - a^\star) u_{i,\sigma}^n u_{j,\sigma}^n (D_{K\sigma} \ln(\mathbf{u}_i) - D_{K\sigma} \ln(\mathbf{u}_j)) \right) \\ - a^\star \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma D_{K\sigma} \mathbf{u}_i^n = 0, \quad \forall K \in \mathcal{F}, \forall i \in \llbracket 1, N \rrbracket. \end{aligned} \quad (19)$$

This reformulation is suitable for proving a discrete entropy - entropy dissipation inequality, which should be seen as a discrete counterpart of Proposition 1.3.

**Proposition 3.7** *Given  $\mathbf{U}^{n-1}$  in  $\mathcal{A}^{\mathcal{F}}$ , any solution  $\mathbf{U}^n \in \mathcal{A}^{\mathcal{F}}$  to (14) satisfies*

$$E_{\mathcal{F}}(\mathbf{U}^n) - E_{\mathcal{F}}(\mathbf{U}^{n-1}) + \Delta t_n \min_{1 \leq i, j \leq N} a_{i,j} \sum_{\sigma \in \mathcal{E}} \sum_{i=1}^N \tau_\sigma u_{i,\sigma}^n (D_{K\sigma} \ln(\mathbf{u}_i^n))^2 \leq 0. \quad (20)$$

In particular,  $E_{\mathcal{F}}(\mathbf{U}^n) \leq E_{\mathcal{F}}(\mathbf{U}^{n-1})$ .

*Proof.* Multiplying equation (19) by  $\Delta t_n \ln(u_{i,K}^n)$  and summing over the cells and species leads to:

$$\begin{aligned} \sum_{K \in \mathcal{F}} \sum_{i=1}^N (u_{i,K}^n \ln(u_{i,K}^n) - u_{i,K}^{n-1} \ln(u_{i,K}^n)) m_K + \Delta t_n a^\star \sum_{\sigma \in \mathcal{E}} \sum_{i=1}^N \tau_\sigma u_{i,\sigma}^n (D_{K\sigma} \ln(\mathbf{u}_i^n))^2 \\ - \Delta t_n \sum_{K \in \mathcal{F}} \sum_{i=1}^N \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \left( \sum_{j=1}^N (a_{i,j} - a^\star) u_{j,\sigma}^n u_{i,\sigma}^n \ln(u_{i,K}^n) D_{K\sigma} (\ln(\mathbf{u}_i^n) - \ln(\mathbf{u}_j^n)) \right) = 0. \end{aligned} \quad (21)$$

Using the symmetry of the matrix  $A$  and discrete integration by part, both in space and with respect to the species, we have:

$$\begin{aligned} \sum_{K \in \mathcal{F}} \sum_{i=1}^N \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \left( \sum_{j=1}^N (a_{i,j} - a^\star) u_{j,\sigma}^n u_{i,\sigma}^n \ln(u_{i,K}^n) D_{K\sigma} (\ln(\mathbf{u}_i^n) - \ln(\mathbf{u}_j^n)) \right) = \\ - \sum_{\sigma \in \mathcal{E}} \tau_\sigma \left( \sum_{1 \leq i < j \leq N} (a_{i,j} - a^\star) u_{j,\sigma}^n u_{i,\sigma}^n (D_{K\sigma} (\ln(\mathbf{u}_i^n) - \ln(\mathbf{u}_j^n)))^2 \right). \end{aligned} \quad (22)$$

On the other hand, the convexity of  $c \ln(c)$  yields:

$$u_{i,K}^n - u_{i,K}^{n-1} + u_{i,K}^n \ln(u_{i,K}^n) - u_{i,K}^{n-1} \ln(u_{i,K}^n) \geq u_{i,K}^n \ln(u_{i,K}^n) - u_{i,K}^{n-1} \ln(u_{i,K}^{n-1}).$$

Combining this inequality with Equation (22) and Lemma 3.5 in (21) provides:

$$\begin{aligned} E_{\mathcal{F}}(\mathbf{U}^n) - E_{\mathcal{F}}(\mathbf{U}^{n-1}) + \Delta t_n a^\star \sum_{\sigma \in \mathcal{E}} \sum_{i=1}^N \tau_\sigma u_{i,\sigma}^n (D_{K\sigma} \ln(\mathbf{u}_i^n))^2 \\ + \Delta t_n \sum_{\sigma \in \mathcal{E}} \tau_\sigma \left( \sum_{1 \leq i < j \leq N} (a_{i,j} - a^\star) u_{j,\sigma}^n u_{i,\sigma}^n (D_{K\sigma} (\ln(\mathbf{u}_i^n) - \ln(\mathbf{u}_j^n)))^2 \right) \leq 0. \end{aligned}$$

Using the hypothesis  $0 \leq \min a_{i,j} \leq a^*$  together with

$$\sum_{i=1}^N u_{i,\sigma}^n (D_{K\sigma} \ln(\mathbf{u}_i^n))^2 - \left( \sum_{1 \leq i < j \leq N} u_{j,\sigma}^n u_{i,\sigma}^n (D_{K\sigma} (\ln(\mathbf{u}_i^n) - \ln(\mathbf{u}_j^n)))^2 \right) = \sum_{i=1}^N u_{i,\sigma}^n \left( 1 - \sum_{j=1}^N u_{j,\sigma}^n \right) (D_{K\sigma} \ln(\mathbf{u}_i^n))^2 \geq 0, \quad (23)$$

we deduce that (20) holds.  $\square$

The proof of Theorem 2.2 is now complete.

## 4 Convergence analysis

The goal of this Section is to prove Theorem 2.3, which states the convergence of the approximate solution towards a weak solution to the continuous problem in the sense of Definition 1.4 under the nondegeneracy condition (7). We could extend this result on several other special cases including the one treated in [7]. We hint that the optimal assumption would be that the zeros of the diffusion matrix form a cluster-graph. However, we stick to the study of the non-degenerate case for the sake of simplicity.

We consider here a sequence  $(\mathcal{T}_m, \mathcal{E}_m, (x_K)_{K \in \mathcal{T}_m})_{m \geq 1}$  of admissible discretizations with  $h_{\mathcal{T}_m}, h_{T,m}$  tending to 0 as  $m$  tends to  $+\infty$ , while the regularity  $\zeta_{\mathcal{T}_m}$  remains uniformly bounded from below by a positive constant  $\zeta^*$ . Theorem 2.2 provides the existence of a family of discrete solutions  $\mathbf{U}_m = (u_{i,K}^n)_{i \in \llbracket 1, N \rrbracket, K \in \mathcal{T}_m, 1 \leq n \leq N_m}$ . To prove Theorem 2.3, we first establish in Section 4.2 some compactness properties on the family of piecewise constant approximate solutions  $U_{\mathcal{T}_m, \Delta t_m}$ . Then we identify the limit as a weak solution in Section 4.3. In order to enlighten the notations, we remove the subscript  $m$  as soon as it is not necessary for understanding.

### 4.1 Reconstruction operators

To carry out the convergence analysis, we introduce some reconstruction operators following the methodology proposed in [16]. The operators  $\pi_{\mathcal{T}} : \mathbb{R}^{\mathcal{T}} \rightarrow L^\infty(\Omega)$  and  $\pi_{\mathcal{T}, \Delta t} : (\mathbb{R}^{\mathcal{T}})^{N_T} \rightarrow L^\infty(Q_T)$  are defined respectively by

$$\pi_{\mathcal{T}} \mathbf{f}(x) = f_K \quad \text{if } x \in K, \quad \forall \mathbf{f} = (f_K)_{K \in \mathcal{T}},$$

and

$$\pi_{\mathcal{T}, \Delta t} \mathbf{f}(t, x) = f_K^n \quad \text{if } (t, x) \in (t_{n-1}, t_n] \times K, \quad \forall \mathbf{f} = (f_K^n)_{K \in \mathcal{T}, 1 \leq n \leq N_T}.$$

These operators allow to pass from the discrete solution  $(\mathbf{U}^n)_{1 \leq n \leq N_T}$  to the approximate solution since

$$u_{i,\mathcal{T}, \Delta t} = \pi_{\mathcal{T}, \Delta t} (\mathbf{u}_i^n)_n, \quad \forall i \in \llbracket 1, N \rrbracket.$$

In order to carry out the analysis, we further need to introduce approximate gradient reconstruction. For  $\sigma = K|L \in \mathcal{E}_{\text{int}}$ , we denote by  $\Delta_\sigma$  the diamond cell corresponding to  $\sigma$ , which is the interior of the convex hull of  $\{\sigma, x_K, x_L\}$ . For  $\sigma \in \mathcal{E}_{\text{ext}}$ , the diamond cell  $\Delta_\sigma$  is defined as the interior of the convex hull of  $\{\sigma, x_K\}$ . The approximate gradient  $\nabla_{\mathcal{T}} : \mathbb{R}^{\mathcal{T}} \rightarrow L^2(\Omega)^d$  we use in the analysis is merely weakly

consistent (unless  $d = 1$ ) and takes its source in [10, 17]. It is piecewise constant on the diamond cells  $\Delta_\sigma$ , and it is defined as follows:

$$\nabla_{\mathcal{F}} \mathbf{f}(x) = d \frac{D_{K\sigma} \mathbf{f}}{d_\sigma} n_{K\sigma} \quad \text{if } x \in \Delta_\sigma, \quad \forall \mathbf{f} \in \mathbb{R}^{\mathcal{F}},$$

where  $n_{K\sigma}$  is the outer-pointing normal of  $K$  at  $\sigma$ . We also define  $\nabla_{\mathcal{F}, \Delta t} : \mathbb{R}^{\mathcal{F} \times N_T} \rightarrow L^2(Q_T)^d$  by setting

$$\nabla_{\mathcal{F}, \Delta t} \mathbf{f}(t, \cdot) = \nabla_{\mathcal{F}} \mathbf{f}^n \quad \text{if } t \in (t_{n-1}, t_n], \quad \forall \mathbf{f} = (\mathbf{f}^n)_{1 \leq n \leq N_T} \in \mathbb{R}^{\mathcal{F} \times N_T}.$$

It follows from the definition of the approximate gradient that

$$\sum_{\sigma \in \mathcal{E}} \tau_\sigma D_{K\sigma} \mathbf{f} D_{K\sigma} \mathbf{g} = \frac{1}{d} \int_{\Omega} \nabla_{\mathcal{F}} \mathbf{f} \cdot \nabla_{\mathcal{F}} \mathbf{g} dx, \quad \forall \mathbf{f}, \mathbf{g} \in \mathbb{R}^{\mathcal{F}}. \quad (24)$$

This implies in particular that

$$\sum_{\sigma \in \mathcal{E}} \tau_\sigma |D_\sigma \mathbf{f}|^2 = \frac{1}{d} \int_{\Omega} |\nabla_{\mathcal{F}} \mathbf{f}|^2 dx, \quad \forall \mathbf{f} \in \mathbb{R}^{\mathcal{F}}. \quad (25)$$

## 4.2 Compactness properties

In this subsection, we take advantage of Proposition 3.7 and of the non-degeneracy assumption (7) to get enough compactness for the convergence.

**Lemma 4.1** *There exists  $C$  depending only on  $\Omega$  and  $\min_{i \neq j} a_{i,j}$  such that*

$$\sum_{i=1}^N \iint_{Q_T} |\nabla_{\mathcal{F}_m, \Delta t_m} \sqrt{\mathbf{u}_{i,m}}|^2 + (\pi_{\mathcal{F}_m, \Delta t_m} \sqrt{\mathbf{u}_{i,m}})^2 dx dt \leq C, \quad \forall m \geq 1.$$

*Proof.* We get rid of the subscript  $m$  for the ease of reading. The  $L^\infty$  bound on  $\mathbf{U}$  yields immediately the  $L^2$  estimate on  $\pi_{\mathcal{F}, \Delta t} \sqrt{\mathbf{u}_i}$ . The proof thus consists in proving the bound on the discrete gradient. Let us focus on the proof of  $\iint_{Q_T} |\nabla_{\mathcal{F}, \Delta t} \sqrt{\mathbf{u}_i}|^2 dx dt \leq C$  for some fixed  $i \in \llbracket 1, N \rrbracket$ . Thanks to (25), we have

$$\begin{aligned} \iint_{Q_T} |\nabla_{\mathcal{F}, \Delta t} \sqrt{\mathbf{u}_i}|^2 &= d \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}_{\text{int}}} \tau_\sigma |D_\sigma \sqrt{\mathbf{u}_i^n}|^2, \\ &= d \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}_{\text{int}}} \tau_\sigma \check{u}_{i\sigma}^n |D_\sigma \ln(\mathbf{u}_i^n)|^2, \end{aligned}$$

where  $\check{u}_{i\sigma}^n = 4 \frac{(D_\sigma \sqrt{\mathbf{u}_i^n})^2}{(D_\sigma \ln(\mathbf{u}_i^n))^2}$ . It results from Cauchy-Schwarz inequality that  $\check{u}_{i\sigma}^n \leq u_{i\sigma}^n$ . Therefore, Proposition 3.7 provides:

$$\min_{i \neq j} a_{i,j} \sum_{i=1}^N \iint_{Q_T} |\nabla_{\mathcal{F}, \Delta t} \sqrt{\mathbf{u}_i}|^2 \leq \frac{d}{4} (E_{\mathcal{F}}(\mathbf{U}^0) - E_{\mathcal{F}}(\mathbf{U}^{N_T})).$$

As  $E_{\mathcal{F}}$  is bounded between  $-m_\Omega$  and 0 and as, by hypothesis,  $\min a_{i,j} > 0$ , we obtain the desired bound.  $\square$

The inequality  $2D_\sigma \sqrt{\mathbf{u}_i^n} \geq D_\sigma \mathbf{u}_i^n$  and Lemma 4.1 yield the following discrete  $L^2(0, T; H^1(\Omega))$  estimate on  $\mathbf{u}_i$ .

**Corollary 4.2** *There exists  $C$  depending only on  $\Omega$  and  $\min_{i \neq j} a_{i,j}$  such that*

$$\sum_{i=1}^N \iint_{Q_T} |\nabla_{\mathcal{T}_m, \Delta t_m} \mathbf{u}_{i,m}|^2 + (\pi_{\mathcal{T}_m, \Delta t_m} \mathbf{u}_{i,m})^2 dx dt \leq C, \quad \forall m \geq 1.$$

The following proposition is about the relative compactness of the approximate solution and of the weakly consistent approximate gradient.

**Proposition 4.3** *Let  $(\mathbf{U}_m)$  be the family of discrete solutions. There exists at least one  $U \in L^\infty(Q_T; \mathcal{A}) \cap L^2((0, T); H^1(\Omega))$  such that, up to a subsequence, for all  $i \in \llbracket 1, N \rrbracket$ :*

$$\pi_{\mathcal{T}_m, \Delta t_m} \mathbf{u}_{i,m} \xrightarrow{m \rightarrow \infty} u_i \quad \text{strongly in } L^2(Q_T), \quad (26)$$

$$\nabla_{\mathcal{T}_m, \Delta t_m} \mathbf{u}_{i,m} \xrightarrow{m \rightarrow \infty} \nabla u_i \quad \text{weakly in } L^2(Q_T)^d. \quad (27)$$

*Proof.* We drop the subscript  $m$  for clarity. The proof of this result relies on a discrete Aubin-Lions lemma [20, Lemma 3.4] on the particular setting of [7, Lemma 9]. Define the discrete  $L^2(0, T; (H^1(\Omega))')$  norm by duality as follows:

$$\|\mathbf{v}\|_{-1} = \sup \left\{ \int_{\Omega} \pi_{\mathcal{T}} \mathbf{v} \pi_{\mathcal{T}} \boldsymbol{\varphi}, \|\pi_{\mathcal{T}} \boldsymbol{\varphi}\|_{L^2}^2 + \|\nabla_{\mathcal{T}} \boldsymbol{\varphi}\|_{L^2}^2 = 1 \right\}, \quad \forall \mathbf{v} \in \mathbb{R}^{\mathcal{T}}.$$

Therefore if  $\|\nabla_{\mathcal{T}, \Delta t} \mathbf{u}_i\|_{L^2(Q_T)} \leq C$  and  $\sum_n \|\mathbf{u}_i^n - \mathbf{u}_i^{n-1}\|_{-1} \leq C$ , then, up to a subsequence,  $\pi_{\mathcal{T}, \Delta t} \mathbf{u}_i$  tends towards some  $u_i$  in  $L^2(Q_T)$ , while  $\nabla_{\mathcal{T}, \Delta t} \mathbf{u}_i$  converges weakly towards  $\nabla u_i$ . In particular,  $U \in L^2(0, T; H^1(\Omega))^N$ .

Corollary 4.2 provides the  $L^2$  bound on  $\nabla_{\mathcal{T}, \Delta t} \mathbf{u}_i$ . For the other inequality, we let  $\boldsymbol{\varphi} \in \mathbb{R}^{\mathcal{T}}$ ,  $n \in \llbracket 1, N_T \rrbracket$  and  $i \in \llbracket 1, N \rrbracket$ . It follows from (14a) that

$$\int_{\Omega} \pi_{\mathcal{T}} (\mathbf{u}_i^n - \mathbf{u}_i^{n-1}) \pi_{\mathcal{T}} \boldsymbol{\varphi} = -\Delta t_n \sum_{K \in \mathcal{T}} \varphi_K \sum_{\sigma \in \mathcal{E}_K} F_{i, K\sigma}^n.$$

Using (14b), this yields

$$\begin{aligned} \frac{1}{\Delta t_n} \int_{\Omega} \pi_{\mathcal{T}} (\mathbf{u}_i^n - \mathbf{u}_i^{n-1}) \pi_{\mathcal{T}} \boldsymbol{\varphi} &= \sum_{\sigma \in \mathcal{E}} a^* \tau_{\sigma} D_{K\sigma} \mathbf{u}_i^n D_{K\sigma} \boldsymbol{\varphi} \\ &\quad + \sum_{\sigma \in \mathcal{E}} \tau_{\sigma} \left( \sum_{j=1}^N (a_{i,j} - a^*) (u_{j,\sigma}^n D_{K\sigma} u_i^n - u_{i,\sigma}^n D_{K\sigma} u_j^n) \right) D_{K\sigma} \boldsymbol{\varphi}. \end{aligned}$$

Using the Cauchy-Schwarz inequality, the  $L^\infty$  bound on  $(u_{i,\sigma}^n)_{\sigma \in \mathcal{E}, i \in \llbracket 1, N \rrbracket}$  and (24) then leads to

$$\begin{aligned} \frac{1}{\Delta t_n} \int_{\Omega} \pi_{\mathcal{T}} (\mathbf{u}_i^n - \mathbf{u}_i^{n-1}) \pi_{\mathcal{T}} \boldsymbol{\varphi} &\leq a^* \|\nabla_{\mathcal{T}} \mathbf{u}_i^n\|_{L^2(\Omega)} \|\nabla_{\mathcal{T}} \boldsymbol{\varphi}\|_{L^2(\Omega)} \\ &\quad + \|\nabla_{\mathcal{T}} \boldsymbol{\varphi}\|_{L^2(\Omega)} \sum_{j=1}^N |a_{i,j} - a^*| \left( \|\nabla_{\mathcal{T}} \mathbf{u}_i^n\|_{L^2(\Omega)} + \|\nabla_{\mathcal{T}} \mathbf{u}_j^n\|_{L^2(\Omega)} \right). \end{aligned}$$

By definition of the discrete  $(H^1(\Omega))'$  norm, we have

$$\left\| \frac{\mathbf{u}_i^n - \mathbf{u}_i^{n-1}}{\Delta t_n} \right\|_{-1} \leq a^* \|\nabla_{\mathcal{T}} \mathbf{u}_i^n\|_{L^2(\Omega)} + \sum_{j=1}^N |a_{i,j} - a^*| \left( \|\nabla_{\mathcal{T}} \mathbf{u}_i^n\|_{L^2(\Omega)} + \|\nabla_{\mathcal{T}} \mathbf{u}_j^n\|_{L^2(\Omega)} \right).$$



Using Corollary 4.2 again provides that  $\sum_n \|\mathbf{u}_i^n - \mathbf{u}_i^{n-1}\|_{-1} \leq C$ . The relative compactness properties on  $\pi_{\mathcal{T}, \Delta t} \mathbf{u}_i$  and  $\nabla_{\mathcal{T}, \Delta t} \mathbf{u}_i$  follow.

We still have to prove that  $U$  is in  $L^\infty(Q_T; \mathcal{A})$ . Let  $i \in \llbracket 1, N \rrbracket$  and let  $\varphi_i \in L^2(Q_T)$  be zero where the limit  $u_i$  is nonnegative and 1 where the limit is negative, then

$$\int_{Q_T} \varphi_i \pi_{\mathcal{T}, \Delta t} \mathbf{u}_i \geq 0 \quad \text{and} \quad \int_{Q_T} \varphi_i \pi_{\mathcal{T}, \Delta t} \mathbf{u}_i \xrightarrow{m \rightarrow +\infty} \int_{Q_T} u_i \varphi_i \leq 0.$$

Therefore,  $\int_{Q_T} u_i \varphi_i = 0$ , so that  $u_i$  is nonnegative. Finally, the linearity of the limit yields  $\sum_{i=1}^N u_i = 1$ .  $\square$

**Remark 4.1** *The uniform  $L^\infty(Q_T)$  bound on  $\pi_{\mathcal{T}_m, \Delta t_m} \mathbf{U}_m$  together with the strong convergence in  $L^2(Q_T)$  yield (16) thanks to Hölder's inequality:*

$$\pi_{\mathcal{T}_m, \Delta t_m} \mathbf{U}_m \xrightarrow{m \rightarrow \infty} U \quad \text{strongly in } L^p(Q_T)^N, \text{ for any } 1 \leq p < \infty.$$

We also need convergence properties for the face values  $u_{i,\sigma}$ . We can reconstruct an approximate solution  $u_{i,\mathcal{E}, \Delta t}$  which is piecewise constant on the diamond cells by setting, for all  $i \in \llbracket 1, N \rrbracket$ :

$$u_{i,\mathcal{E}, \Delta t}(t, x) = u_{i,\sigma}^n \quad \text{if } (t, x) \in (t_{n-1}, t_n] \times \Delta_\sigma, \quad \sigma \in \mathcal{E}.$$

**Lemma 4.4** *We have, for any  $i \in \llbracket 1, N \rrbracket$ :*

$$u_{i,\mathcal{E}_m, \Delta t_m} \xrightarrow{m \rightarrow \infty} u_i \quad \text{in } L^p(Q_T), \text{ for any } 1 \leq p < \infty,$$

where  $U$  is as in Proposition 4.3.

*Proof.* Here again, we get rid of  $m$  for clarity, and show the convergence for a specific value of  $p$ . The convergence for any finite  $p$  follows from the  $L^\infty(Q_T)$  bound on  $u_{i,\mathcal{E}_m, \Delta t_m}$  and Hölder's inequality. Since  $u_{i,\mathcal{T}, \Delta t}$  converges towards  $u_i$  in  $L^1(Q_T)$ , and since  $u_{i,\mathcal{E}, \Delta t}$  is uniformly bounded, it suffices to show that  $\|u_{i,\mathcal{E}, \Delta t} - u_{i,\mathcal{T}, \Delta t}\|_{L^1(Q_T)}$  tends to 0. Denote by  $\Delta_{K\sigma}$  the half-diamond cell which is defined as the interior of the convex hull of  $\{x_K, \sigma\}$  for  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$ , then the following geometrical relation holds:

$$m(\Delta_{K\sigma}) = \frac{1}{d} m_\sigma \text{dist}(x_K, \sigma) \leq \frac{h_{\mathcal{T}}}{d} m_\sigma.$$

As a consequence,

$$\begin{aligned} \|u_{i,\mathcal{E}, \Delta t} - u_{i,\mathcal{T}, \Delta t}\|_{L^1(Q_T)} &= \sum_{n=1}^{N_T} \Delta t_n \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_{\Delta_{K\sigma}} |u_{i,K}^n - u_{i,\sigma}^n| \\ &\leq \frac{h_{\mathcal{T}}}{d} \sum_{n=1}^{N_T} \Delta t_n \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma |u_{i,K}^n - u_{i,\sigma}^n|. \end{aligned}$$

As we have  $u_{i,K}^n = u_{i,\sigma}^n$ , the contributions corresponding to the boundary edges vanish. For  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $u_{i,\sigma}$  is an average of  $u_{i,K}$  and  $u_{i,K\sigma}$ , hence  $|u_{i,K}^n - u_{i,\sigma}^n| \leq |u_{i,K}^n - u_{i,K\sigma}^n|$ . Therefore, we obtain that

$$\begin{aligned} \|u_{i,\mathcal{E}, \Delta t} - u_{i,\mathcal{T}, \Delta t}\|_{L^1(Q_T)} &\leq \frac{h_{\mathcal{T}}}{d} \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}_{\text{int}}} 2m_\sigma |D_\sigma \mathbf{u}_i^n| \\ &\leq 2 \frac{h_{\mathcal{T}}}{d} \left( \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}_{\text{int}}} m_\sigma d_\sigma \right)^{\frac{1}{2}} \left( \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}_{\text{int}}} \tau_\sigma |D_\sigma \mathbf{u}_i^n|^2 \right)^{\frac{1}{2}}. \end{aligned}$$

We deduce from Corollary 4.2 that  $\|u_{i,\mathcal{E}, \Delta t} - u_{i,\mathcal{T}, \Delta t}\|_{L^1(Q_T)} \leq Ch_{\mathcal{T}}$ , hence  $u_{i,\mathcal{E}, \Delta t}$  and  $u_{i,\mathcal{T}, \Delta t}$  share the same limit in  $L^1(Q_T)$ .  $\square$

### 4.3 Convergence towards a weak solution

The last step to conclude the proof of Theorem 2.3 is to identify the limit value  $U$  exhibited in Proposition 4.3 as a weak solution to (1), (3) corresponding to the initial profile  $U \in L^\infty(\Omega; \mathcal{A})$ . This is the purpose of our last statement.

**Proposition 4.5** *Let  $U$  be as in Proposition 4.3, then  $U$  is a weak solution in the sense of Definition 1.4.*

*Proof.* We drop again the subscript  $m$  for the sake of readability, and let  $i \in \llbracket 1, N \rrbracket$ ,  $\varphi \in C_c^\infty([0, T] \times \overline{\Omega})$ , then define  $\boldsymbol{\varphi} = (\varphi_K^n)$  by  $\varphi_K^n = \varphi(x_K, t_n)$  for all  $n \in \{0, \dots, N_T\}$  and  $K \in \mathcal{T}$ . Multiplying (14a) by  $\Delta t_n \varphi_K^{n-1}$ , then summing over  $K \in \mathcal{T}$  and  $n \in \{1, \dots, N_T\}$  leads to

$$T_1 + T_2 + T_3 = 0, \quad (28)$$

where we have set

$$\begin{aligned} T_1 &= \sum_{n=1}^{N_T} \sum_{K \in \mathcal{T}} m_K (u_{i,K}^n - u_{i,K}^{n-1}) \varphi_K^{n-1}, \\ T_2 &= \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}} \tau_\sigma a^\star D_{K\sigma} \mathbf{u}_i^n D_{K\sigma} \boldsymbol{\varphi}^{n-1}, \\ T_3 &= \sum_{n=1}^{N_T} \Delta t_n \sum_{\sigma \in \mathcal{E}} \tau_\sigma \sum_{j=1}^N (a_{i,j} - a^\star) \left( u_{j\sigma}^n D_{K\sigma} \mathbf{u}_i^n - u_{i\sigma}^n D_{K\sigma} \mathbf{u}_j^n \right) D_{K\sigma} \boldsymbol{\varphi}^{n-1}. \end{aligned}$$

The term  $T_1$  can be rewritten as

$$T_1 = \sum_{n=1}^{N_T} \Delta t_n \sum_{K \in \mathcal{T}} m_K u_{i,K}^n \frac{\varphi_K^{n-1} - \varphi_K^n}{\Delta t_n} - \sum_{K \in \mathcal{T}} m_K u_{i,K}^0 \varphi_K^0,$$

so that it follows from the convergence of  $\pi_{\mathcal{T}, \Delta t} \mathbf{U}$  towards  $U$  and of  $\pi_{\mathcal{T}} \mathbf{U}^0$  towards  $U^0$  together with the regularity of  $\varphi$  that

$$T_1 \xrightarrow{m \rightarrow \infty} - \iint_{Q_T} u_i \partial_t \varphi dx dt - \int_{\Omega} u_i^0 \varphi(0, \cdot) dx. \quad (29)$$

To treat the term  $T_2$ , we introduce a strongly consistent reconstruction of the gradient. Following [15] (see [11] for a practical example), one can reconstruct a second approximate gradient operator  $\widehat{\nabla}_{\mathcal{T}} : \mathbb{R}^{\mathcal{T}} \rightarrow L^\infty(\Omega)^d$  such that

$$\int_{\Delta_\sigma} \nabla_{\mathcal{T}} \mathbf{u} \cdot \widehat{\nabla}_{\mathcal{T}} \mathbf{v} dx = \tau_\sigma D_{K\sigma} \mathbf{u} D_{K\sigma} \mathbf{v}, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^{\mathcal{T}}, \forall \sigma \in \mathcal{E},$$

and which is strongly consistent, i.e.,

$$\widehat{\nabla}_{\mathcal{T}} \boldsymbol{\varphi}^n \xrightarrow{h_{\mathcal{T}} \rightarrow 0} \nabla \varphi(\cdot, t_n) \text{ uniformly in } \overline{\Omega}, \quad \forall n \in \{1, \dots, N_T\},$$

thanks to the smoothness of  $\varphi$ . Using this tool, the terms  $T_2$  and  $T_3$ , are easy to treat. The first one can be rewritten as:

$$T_2 = a^\star \iint_{Q_T} \nabla_{\mathcal{T}, \Delta t} \mathbf{u}_i \cdot \widehat{\nabla}_{\mathcal{T}, \Delta t} \boldsymbol{\varphi} dx dt,$$

so that

$$T_2 \xrightarrow{m \rightarrow \infty} a^* \iint_{Q_T} \nabla u_i \cdot \nabla \varphi \, dx \, dt. \quad (30)$$

On the other hand, the term  $T_3$  rewrites

$$T_3 = \iint_{Q_T} \sum_{j=1}^N (a_{i,j} - a^*) (u_{j,\mathcal{E},\Delta t} \nabla_{\mathcal{F},\Delta t} \mathbf{u}_i - u_{i,\mathcal{E},\Delta t} \nabla_{\mathcal{F},\Delta t} \mathbf{u}_j) \widehat{\nabla}_{\mathcal{F},\Delta t} \varphi,$$

so that

$$T_3 \xrightarrow{m \rightarrow \infty} \iint_{Q_T} \sum_{j=1}^N (a_{i,j} - a^*) (u_j \nabla u_i - u_i \nabla u_j) \nabla \varphi. \quad (31)$$

Combining (28), (29), (30), and (31), we obtain that

$$\begin{aligned} & - \iint_{Q_T} u_i \partial_t \varphi \, dx \, dt - \int_{\Omega} u_i^0 \varphi(0, \cdot) \, dx + a^* \iint_{Q_T} \nabla u_i \cdot \nabla \varphi \, dx \, dt \\ & \quad + \iint_{Q_T} \sum_{j=1}^N (a_{i,j} - a^*) (u_j \nabla u_i - u_i \nabla u_j) \nabla \varphi \, dx \, dt = 0, \quad \forall \varphi \in C_c^\infty([0, T] \times \overline{\Omega}). \end{aligned}$$

Using  $U \in \mathcal{A}$  and the relation (9), we recover the weak formulation (8).  $\square$

## 5 Numerical results

The numerical scheme has been implemented using MATLAB. The nonlinear system corresponding to the scheme is solved thanks to a variation of the Newton method with stopping criterion  $\|U^{n,k+1} - U^{n,k}\|_\infty \leq 10^{-12}$ . The solution of the Newton iteration,  $U^{n,k+1/3}$ , is then “projected” on  $\mathcal{A}$  by setting  $U^{n,k+2/3} = \max(U^{n,k+1/3}, 10^{-10}\tau)$ , and then for all  $K \in \mathcal{F}$ :  $U_K^{n,k+1} = U_K^{n,k+2/3} / (\sum_{i=1}^N u_{i,K}^{n,k+2/3})$ .

For the first time step, we also make use of a continuation method based on the intermediate diffusion coefficients  $a_{i,j}^\lambda = \lambda a_{i,j} + (1-\lambda)a^*$  with  $\lambda \in [0, 1]$ . The parameter  $\lambda$  is originally set to 1. If the Newton’s method does not converge, we let  $\lambda = (\lambda + \lambda_{\text{prev}})/2$  where  $\lambda_{\text{prev}}$  is originally set to 0. If the Newton’s method converges, we let  $\lambda_{\text{prev}} = \lambda$  and  $\lambda = 1$ .

### 5.1 Convergence under grid refinement

Our first test case is devoted to the convergence analysis of the scheme in a one-dimensional setting  $\Omega = (0, 1)$ . Two different initial conditions are considered:  $U_s^0$  is smooth and vanished point-wise at the boundary of  $\Omega$ , whereas  $U_r^0$  is discontinuous and vanishes on intervals of  $\Omega$ :

$$\begin{aligned} u_{1,s}^0(x) &= \frac{1}{4} + \frac{1}{4} \cos(\pi x), & u_{2,s}^0(x) &= \frac{1}{4} + \frac{1}{4} \cos(\pi x), & u_{3,s}^0(x) &= \frac{1}{2} - \frac{1}{2} \cos(\pi x), \\ u_{1,r}^0 &= \mathbf{1}_{[\frac{3}{8}, \frac{5}{8}]}, & u_{2,r}^0 &= \mathbf{1}_{(\frac{1}{8}, \frac{3}{8})} + \mathbf{1}_{(\frac{5}{8}, \frac{7}{8})}, & u_{3,r}^0 &= \mathbf{1}_{[0, \frac{1}{8}]} + \mathbf{1}_{[\frac{7}{8}, 1]}. \end{aligned}$$

We also consider two cross-diffusion coefficients matrices, one called regular with positive off-diagonal coefficients and another called singular with a few null off-diagonal coefficients:

$$A^{\text{reg}} = \begin{pmatrix} 0 & 0.2 & 1 \\ 0.2 & 0 & 0.1 \\ 1 & 0.1 & 0 \end{pmatrix}, \quad A^{\text{sing}} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 0.1 \\ 1 & 0.1 & 0 \end{pmatrix}.$$

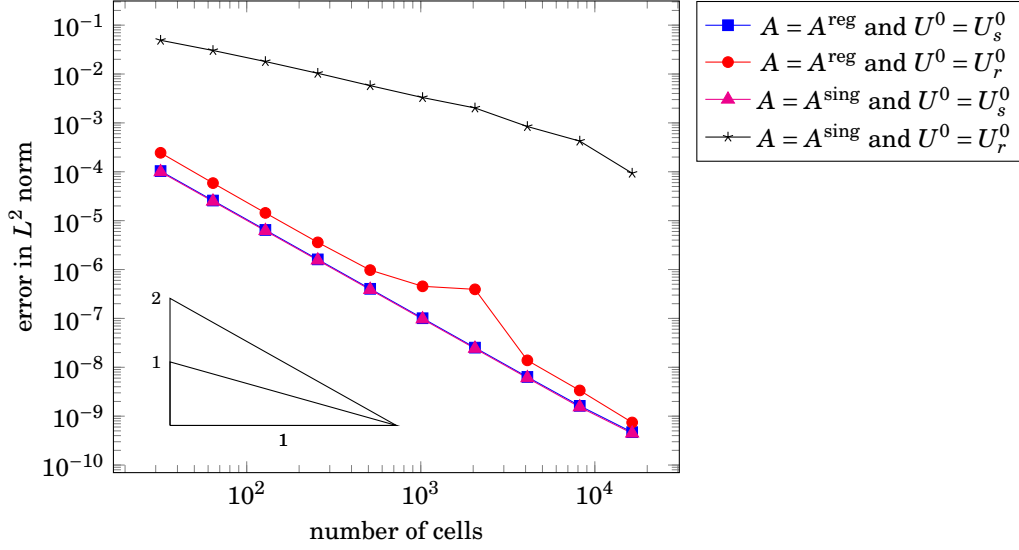


Figure 2: Error with respect to the solution computed on the finest mesh for 1D settings.

For the convergence tests, we have let  $a^* = 0.1$  and the meshes are uniform discretisations of  $[0, 1]$  from  $2^5$  cells to  $2^{15}$  cells. Since we do not have an analytical solution at hand, the approximate solutions are compared to a reference solution computed on the finest grid ( $2^{15}$  cells). The final time is 0.25, and the time discretisation is fixed with a time step of  $2^{-18}$ . Results are summarized in Figure 2. One notices that our scheme is second-order accurate in the setting presented in this paper ( $A = A^{\text{reg}}$ ), but only first-order accurate when confronted to non-diffusive discontinuities. We call non-diffusive discontinuities a spatial discontinuity of  $u_1^0$  and  $u_2^0$  (recall that  $a_{1,2} = 0$  in  $A^{\text{sing}}$ ) for which  $u_3^0$  is equal to 0 on both sides of the discontinuity, so that the contributions corresponding to  $a_{1,3}$  and  $a_{2,3}$  vanish at  $t = 0$ . The origin of this lower order may lie in the difficulty to compute accurately the near-zero concentrations in the neighborhood of such discontinuities.

## 5.2 On the influence of the parameter $a^*$

The choice of  $a^*$  is a natural question concerning our scheme. The equation (13) gives a lower bound:  $a^* > 0$ . The existence of an upper bound is not as clear. Equation (23) shows that for large  $a^*$ , we over-estimate the diffusion. The optimal value of  $a^*$  depends on many variables such as the initial condition, the final time, and the mesh. Optimal choices of  $a^*$  are reported in Table 1. Notice that the optimal value is test cases dependent, since it is affected by the initial condition and by the final time.

One notices on Fig. 3 that the dependency of the quality of the results is much stronger for the initial data  $U_r^0$ . This is due to the presence of vanishing concentrations in some cells, so that the choice  $a^* = 0$  would allow for spurious solutions as highlighted in Remark 2.1. In this situation, the choice of  $a^*$  strongly affects the quality of the results, especially for the first time steps where some concentrations are still close to 0. The numerical experiment and homogeneity considerations suggest

			$A = A^{\text{reg}}$		$A = A^{\text{sing}}$	
			$U^0 = U_s^0$	$U^0 = U_r^0$	$U^0 = U_s^0$	$U^0 = U_r^0$
nb. of cells	32	$T = 0.125$	0.86	0.21	0.79	0.0023
		$T = 0.25$	0.67	0.13	0.49	0.00082
	128	$T = 0.125$	0.86	0.17	0.79	0.00050
		$T = 0.25$	0.67	0.11	0.49	0.00049

Table 1: Values of  $a_{\text{opt}}^*$  for different parameters.  $a_{\text{opt}}^*$  is computed with respect to the reference solution of Section 5.1 for the  $L^2$  norm.

the following suboptimal rule for choosing  $a^*$ :

$$a^* = \min \left\{ \max_{i \neq j} a_{i,j}; \max \left\{ \min_{i \neq j} a_{i,j}, \epsilon \frac{h_{\mathcal{G}}^2}{\tau} \right\} \right\},$$

where  $h_{\mathcal{G}}$  is the mesh size,  $\tau$  the current time step and  $\epsilon$  a small parameter to be tuned by the user.

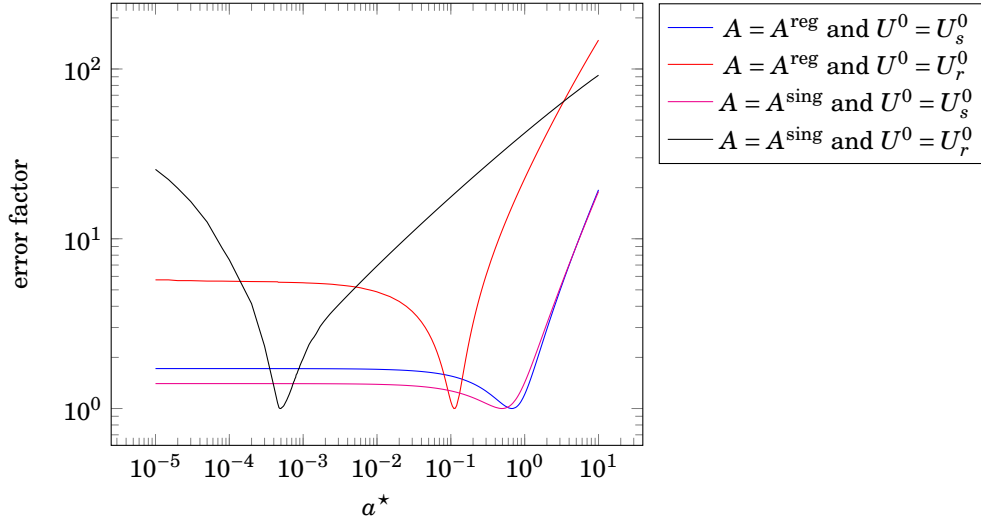


Figure 3: Evolution of the ratio  $\frac{\|U_{a^*} - U_{\text{ref}}\|_2}{\|U_{\text{opt}} - U_{\text{ref}}\|_2}$ , where  $U_{a^*}$  is computed with  $2^7$  cells and  $U_{\text{ref}}$  is as in Section 5.1.

### 5.3 A 2D test case

Our second test is two-dimensional. We choose  $A^{\text{sing}}$  as the diffusion matrix and  $a^* = 0.1$ . The domain  $\Omega = (0, 22) \times (0, 16)$  is discretized into a cartesian grid made of  $110 \times 80$  cells. We use a uniform time stepping with  $\tau = 2^{-3}$ . The initial condition  $U^0$  is depicted in Figure 4. The corresponding steady-state and long-time limit  $U^\infty$  is constant w.r.t. space, i.e.,  $u_i^\infty(x) = \oint u_i^0(y) dy$  for all  $x \in \Omega$ . The time evolution of the relative energy  $E_{\mathcal{G}}(U) - T_{\mathcal{G}}(U^\infty)$  is plotted on Figure 5, showing exponential decay to the

steady-state even though the diffusion matrix is singular. Snapshots showing the evolution of the concentration profiles are presented in Figure 6.

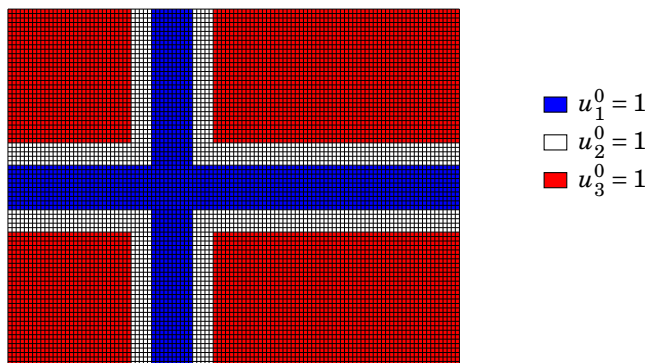


Figure 4: Initial configuration  $U^0$  for the concentrations

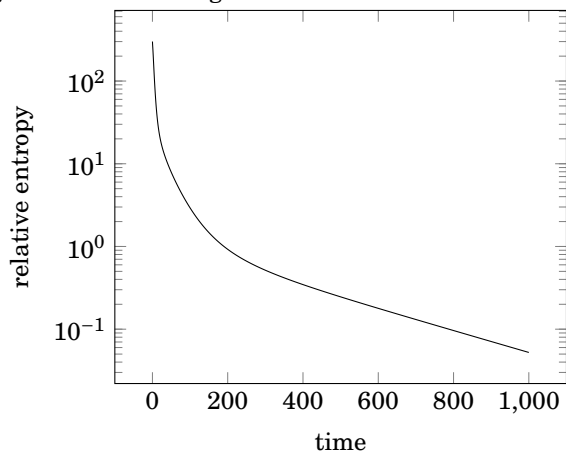


Figure 5:  $E_{\mathcal{F}}(U) - T_{\mathcal{F}}(U^\infty)$  as a function of time.

## 6 Conclusion

We proposed a finite volume scheme based on two-point flux approximation for a degenerate cross-diffusion system. The scheme was designed to preserve the key properties of the continuous system, namely the positivity of the solutions, the constraint on the composition and the decay of the entropy. The scheme requires the introduction of a positive parameter  $a^*$  to avoid unphysical solutions. This parameter plays an important role in the convergence proof, which is carried out under a non-degeneracy assumption. Its importance is also confirmed in the numerical experiments, in particular in the presence of initial profiles with concentrations vanishing in some parts of the computational domain.

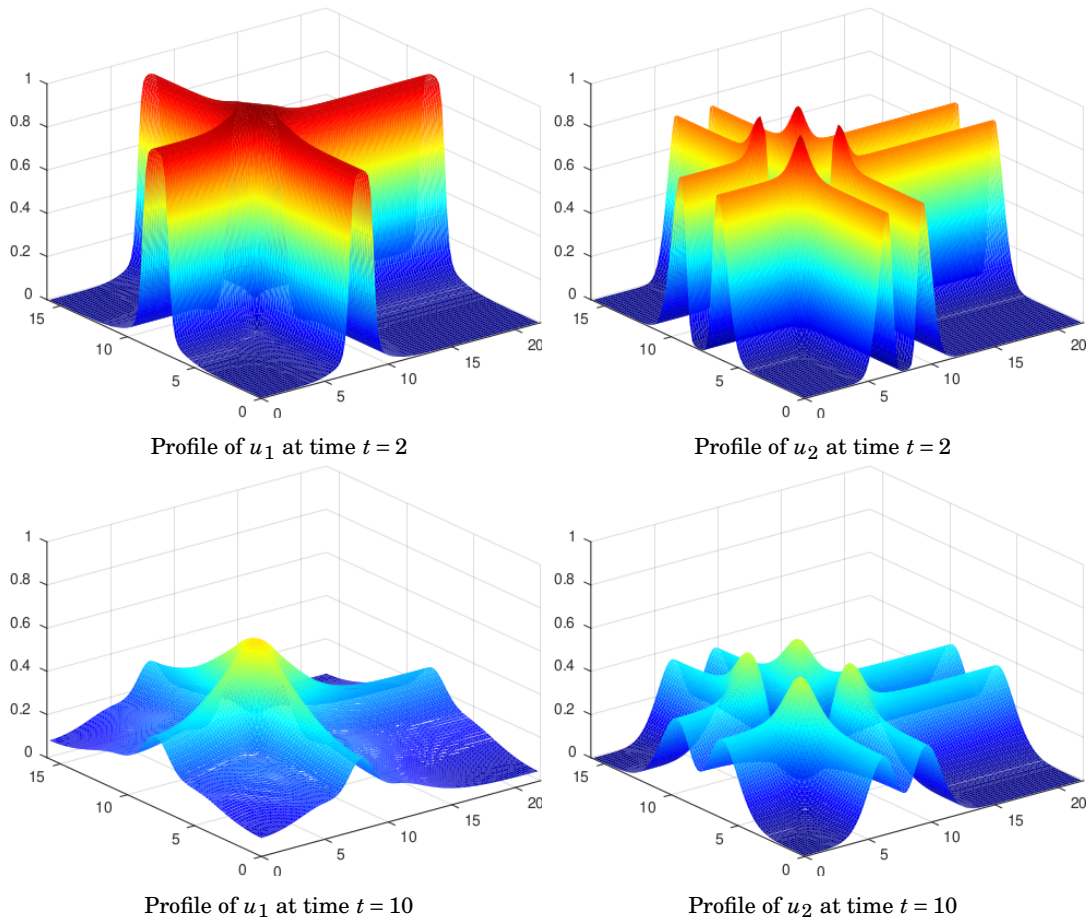


Figure 6: Concentration configurations for various times. The concentration of the third specie can be deduced thanks to  $u_1 + u_2 + u_3 = 1$

## Acknowledgements

The authors acknowledge support from the Labex CEMPI (ANR-11-LABX-0007-01). Clément Cancès also acknowledge support from the COMODO project (ANR-19-CE46-0002), and he warmly thanks Virginie Ehlacher and Laurent Monasse for stimulating discussions that were at the origin of this work.

## 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] A. Ait Hammou Oulhaj and D. Maltese. Convergence of a positive nonlinear control volume finite element scheme for an anisotropic seawater intrusion model with sharp interfaces. *Numer. Methods Partial Differential Equations*, 36(1):133–153, 2019.
- [3] B. Andreianov, M. Bendahmane, and R. Ruiz-Baier. Analysis of a finite volume method for a cross-diffusion model in population dynamics. *Math. Models Methods Appl. Sci.*, 21(2):307–344, 2011.
- [4] A. Bakhta and V. Ehrlacher. Cross-diffusion systems with non-zero flux and moving boundary conditions. *ESAIM: Mathematical Modelling and Numerical Analysis*, 52(4):1385–1415, July 2018.
- [5] J. Berendsen, M. Burger, V. Ehrlacher, and J-F. Pietschmann. Strong solutions and weak-strong stability in a system of cross-diffusion equations. arXiv:1812.10711, 2019.
- [6] M. Burger, B. Schlake, and M.-T. Wolfram. Nonlinear Poisson–Nernst–Planck equations for ion flux through confined geometries. *Nonlinearity*, 25(4):961–990, March 2012.
- [7] C. Cancès, C. Chainais-Hillairet, A. Gerstenmayer, and A. Jüngel. Finite-volume scheme for a degenerate cross-diffusion model motivated from ion transport. *Numerical Methods for Partial Differential Equations*, 35(2):545–575, 2019.
- [8] J. A. Carrillo, F. Filbet, and M. Schmidtchen. Convergence of a finite volume scheme for a system of interacting species with cross-diffusion. arXiv:1804.04385, 2018.
- [9] C. Chainais-Hillairet. Entropy method and asymptotic behaviours of finite volume schemes. In *Finite volumes for complex applications. VII. Methods and theoretical aspects*, volume 77 of *Springer Proc. Math. Stat.*, pages 17–35. Springer, Cham, 2014.
- [10] C. Chainais-Hillairet, J.-G. Liu, and Y.-J. Peng. Finite volume scheme for multi-dimensional drift-diffusion equations and convergence analysis. *ESAIM: M2AN*, 37(2):319–338, 2003.
- [11] Y. Coudière, J.-P. Vila, and P. Villedieu. Convergence rate of a finite volume scheme for a two dimensional convection-diffusion problem. *ESAIM Math. Model. Numer. Anal.*, 33(3):493–516, 1999.
- [12] E. S. Daus, A. Jüngel, and Zurek. Convergence of a finite-volume scheme for a degenerate-singular cross-diffusion system for biofilms. arXiv:2001.09544, 2020.
- [13] K. Deimling. *Nonlinear functional analysis*. Springer-Verlag, Berlin, 1985.
- [14] J. Droniou. Finite volume schemes for diffusion equations: introduction to and review of modern methods. *Math. Models Methods Appl. Sci.*, 24(8):1575–1620, 2014.
- [15] J. Droniou and R. Eymard. The asymmetric gradient discretisation method. In C. Cancès and P. Omnes, editors, *Finite volumes for complex applications VIII - methods and theoretical aspects*, volume 199 of *Springer Proc. Math. Stat.*, pages 311–319, Cham, 2017. Springer.
- [16] J. Droniou, R. Eymard, T. Gallouët, C. Guichard, and R. Herbin. *The Gradient Discretisation Method*, volume 42 of *Mathématiques et Applications*. Springer International Publishing, 2018.



- [17] R. Eymard and T. Gallouët.  $H$ -convergence and numerical schemes for elliptic problems. *SIAM J. Numer. Anal.*, 41(2):539–562, 2003.
- [18] 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.
- [19] 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.
- [20] T. Gallouët and J.-C. Latché. Compactness of discrete approximate solutions to parabolic PDEs—application to a turbulence model. *Commun. Pure Appl. Anal.*, 11(6):2371–2391, 2012.
- [21] 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.
- [22] A. Gerstenmayer and A. Jüngel. Comparison of a finite-element and finite-volume scheme for a degenerate cross-diffusion system for ion transport. arXiv:1812.05849.
- [23] A. Gerstenmayer and A. Jüngel. Analysis of a degenerate parabolic cross-diffusion system for ion transport. *J. Math. Anal. Appl.*, 461(1):523–543, 2018.
- [24] R. Herbin. An error estimate for a finite volume scheme for a diffusion-convection problem on a triangular mesh. *Numer. Methods Partial Differential Equations*, 11(2):165–173, 1995.
- [25] A. Jüngel. The boundedness-by-entropy method for cross-diffusion systems. *Nonlinearity*, 28(6):1963–2001, 2015.
- [26] A. Jüngel. *Entropy methods for diffusive partial differential equations*. SpringerBriefs in Mathematics. Springer, [Cham], 2016.
- [27] A. Jüngel and O. Leingang. Convergence of an implicit Euler Galerkin scheme for Poisson-Maxwell-Stefan systems. *Adv. Comput. Math.*, 45(3):1469–1498, 2019.
- [28] J. Leray and J. Schauder. Topologie et équations fonctionnelles. *Ann. Sci. École Norm. Sup.*, 51((3)):45–78, 1934.
- [29] Z. Sun, J. A. Carrillo, and C.-W. Shu. An entropy stable high-order discontinuous Galerkin method for cross-diffusion gradient flow systems. *Kinet. Relat. Models*, 12(4):885–908, 2019.