**TECHNISCHE UNIVERSITÄT WIEN**

# Dissertation

# Analysis and Numerical Approximation of Degenerate Parabolic Systems arising in Thermodynamics and Biology

ausgeführt zum Zwecke der Erlangung des akademischen Grades

eines Doktors der technischen Wissenschaften unter der Leitung von

## Prof. Dr. Ansgar Jüngel

E101 – Institut für Analysis und Scientific Computing, TU Wien

eingereicht an der Technischen Universität Wien

Fakultät für Mathematik und Geoinformation

von

## Christoph Helmer

Matrikelnummer: 11946209

Diese Dissertation haben begutachtet:

1. **Prof. Dr. Hermann J. Eberl**
   Department of Mathematics and Statistics, University of Guelph

2. **Prof. Dr. Stefanie Sonner**
   Department of Mathematics, Radboud University

Wien, am 11. Mai 2023

# Kurzfassung

Da die Probleme bei der Betrachtung von degenerierten, parabolischen Systemen sehr individuell sind und es keine allgemeine Lösungstheorie gibt, betrachten wir motiviert aus der Anwendung drei verschiedene Modelle.

Das erste Modell kommt aus der Thermodynamik und beschreibt die zeitliche Evolution von Fluiden mit mehreren Komponenten. Wir erweitern die bestehende Literatur um ein Modell, in dem die Temperatur zeit– und ortsabhängig ist und beweisen die globale Existenz schwacher Lösungen unter Ausnutzung der Entropiestruktur, die wir aus der thermodynamischen Modellierung erhalten.

Das zweite Modell kommt aus der Biologie und beschreibt die Evolution von Biofilmen. Wir entwickeln ein Finite–Volumen Schema, für das wir die Existenz und, mit einer zusätzlichen Voraussetzung, die Eindeutigkeit von diskreten Lösungen zeigen. Die Hauptschwierigkeit hierbei besteht in einem degenerierten–singulären Diffusionsterm und dem Beweis der oberen/unteren Schranke für die Biomasse für den wir, anders als im kontinuierlichen Fall, kein Vergleichsprinzip im Diskreten anwenden konnten. Dieses Problem umgehen wir durch die Einführung einer Entropievariablen, welche die gewünschten Schranken garantiert. Des Weiteren zeigen wir, dass die diskreten Lösungen bei der Verfeinerung des Gitters gegen eine schwache Lösung des Systems konvergieren.

Zuletzt betrachten wir ein weiteres Modell, welches die Evolution von Biofilmen beschreibt. Das Modell besteht aus einer degenerierten Reaktions–Diffusionsgleichung und einer lokalen Cahn–Hilliard–Gleichung vierter Ordnung mit degenerierter Mobilität, singulärem Potential und nichtlinearen Quelltermen.

Wir zeigen die globale Existenz schwacher Lösung mit Hilfe einer geeigneten Regularisierung und einer Galerkin–Approximation. Da wir aufgrund der Degeneriertheit keine optimalen Abschätzungen erhalten, um den Limes für die Deregularisierung durchzuführen, benötigen wir einen Minty–Browder–Trick zur Identifikation des Quelltermes.

# Abstract

Since degenerate parabolic systems are quite peculiar, there is no general theory available to obtain the existence of solutions. Thus, we take a closer look at three different models which are motivated applications.

The first model comes from thermodynamics and describes the evolution of multicomponent fluids. We extend the literature by proposing a model which includes nonisothermal temperature as well as Soret/Dufour effects, and prove the global existence by using the entropy structure of the system.

The second model is derived from biology and describes the development of biofilms. We develop a finite–volume scheme for which we prove the existence of discrete solutions and, under additional assumptions, the uniqueness. The main difficulty comes from the degenerate–singular diffusion term and the proof of lower/upper bounds for the biomass fraction since we cannot apply a comparison principle as in the continuous case. We overcome this challenge by introducing an entropy variable which guarantees these bounds. Furthermore, we prove that discrete solutions converge towards a weak solution under mesh refinement.

The last model we discuss is obtained from biology as well and models the growth of biofilms by considering the biomass/solvent as fluid mixture. This system consists of a degenerate reaction–diffusion equation and a local fourth order Cahn–Hilliard equation with degenerate mobility, singular potential and nonlinear source terms. We prove global existence by applying a suitable truncation and a galerkin approximation. Since we do not find optimal estimates due to the degeneracy of the mobility, we perform a Browder–Minty trick for the identification of the source term in the deregularization limit.

# Acknowledgement

First and foremost, I am very grateful to my supervisor, Prof. Dr. Ansgar Jüngel, for giving me the opportunity to join his research group in Vienna. The completion of this thesis would not have been possible without his expertise, patience, and guidance. Despite the challenging circumstances throughout the pandemic, I enjoyed the pleasant and positive work atmosphere he created in the group.

In addition, I would also like to thank Antoine Zurek for the fruitful research collaboration and guidance through the world of finite–volume analysis.

Moving to another country and starting a new job and life has its challenges. Thus, I'm more than thankful to all my former and current colleagues for giving me a warm welcome to Vienna. I also thank all the current members of our group for the very entertaining lunch and non-mathematical discussions which brightened up the atmosphere of our common room.

A special note of thanks goes to Stefan, who fought his way through my Matlab–Code for the Cahn–Hilliard system to "proofread" the Jacobian, to Stefanos for pointing out a mistake in the analysis of the Maxwell–Stefan–Fourier system, to Alexandra, for the brightening discussions throughout Corona–Lockdowns and to Claudia, for the entertaining discussions about basically everything.

I would also like to thank my mother and my sister for their enduring support and encouragement. Last but not least, I am very grateful for the support of my wife, who encourages me to chase my dreams and always believes in me.

# Eidesstattliche Erklärung

Ich erkläre an Eides statt, dass ich die vorliegende Dissertation selbstständig und ohne fremde Hilfe verfasst, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt bzw. die wörtlich oder sinngemäß entnommenen Stellen als solche kenntlich gemacht habe.

Wien, am 11. Mai 2023

_____
Christoph Helmer

# Contents

# Contents

# List of Figures

# List of Tables

# 1 Introduction

The purpose of this thesis is to establish new results in the broad area of analysis/numerical analysis for degenerate parabolic systems. In this context, the term degenerate refers to the diffusion coefficients/the diffusion matrix in the sense, that they are not necessarily strictly positive/positive definite. Although the degeneracy is physically/biologically beneficial in a wide range of contexts, it can cause potential difficulties from a mathematical point of view.

As there is no general theory for degenerate parabolic systems, we have to combine available techniques to individually prove existence for each system. Furthermore, there is no general maximum principle. This would be advantageous as we aspire for naturally lower and upper bounds for the solutions due to the physical/biological context.

In the following chapters, we discuss several degenerate parabolic systems coming from a physical or biological application and present different techniques to handle the lack of estimates due to the degeneracy as well as the proof of upper/lower bounds for the weak solutions or discrete solutions in the finite–volume context. The results in this thesis are based on the publication [HJ21] (Christoph Helmer, Ansgar Jüngel), the publication [HJZ23] (Christoph Helmer, Ansgar Jüngel, Antoine Zurek) and the paper [HJ23](Christoph Helmer, Ansgar Jüngel), which is submitted for publication.

## 1.1 The Maxwell–Stefan–Fourier System

Our goal is to derive and discuss a model, which describes the behavior of a multiple component system in a nonisothermal (no constant temperature) setting, that includes the consideration of Soret and Dufour effects. Mathematically, multicomponent systems are described by Maxwell–Stefan systems and the interaction of the single components with each other is called cross–diffusion. The Soret effect (also known as thermopheresis or thermodiffusion) refers to the particle movement due to a temperature gradient [Lud56]. The Dufour effect is the reciprocal phenomena to the Soret effect and is described by a heat flux of a chemical potential gradient [ME80]. The discussion of these systems have their origin in the 19th century, when they were first described independently by Maxwell for gases [Max67] and Stefan for fluids [Ste71]. For a more detailed overview about physical derivations and applications regarding the Maxwell–Stefan system, we refer to the well–known work of Taylor and Krishna [TK93].

While the physical description goes back to the 19th century, the mathematical existence analysis started only hundred years later in 1998 with [GM98].

Even more recent is the mathematical research for nonisothermal systems, which has been discussed in publications only from 2015 onwards in [GPZ15, PP17, HS18]. However, none of these papers consider the inclusion of the Soret/Dufour effects.

### 1.1.1 The Model Equations

We consider the partial mass densities $\rho_i$ for $i = 1, \ldots, n$ and the temperature $\theta$ in a fluid mixture. This evolution of the system is described as follows:

$$\partial_t \rho_i + \operatorname{div} J_i = r_i, \quad J_i = -\sum_{j=1}^{n} M_{ij}(\boldsymbol{\rho}, \theta) \nabla q_j - M_i(\boldsymbol{\rho}, \theta) \nabla \frac{1}{\theta}, \tag{1.1}$$

$$\partial_t (\rho \theta) + \operatorname{div} J_e = 0, \quad J_e = -\kappa(\theta) \nabla \theta - \sum_{j=1}^{n} M_j(\boldsymbol{\rho}, \theta) \nabla q_j \quad \text{in } \Omega, \; i = 1, \ldots, n, \tag{1.2}$$

where $\Omega \subset \mathbb{R}^3$ is a bounded domain, $\boldsymbol{\rho} = (\rho_1, \ldots, \rho_n)$ is the vector of mass densities, and $q_i = \log(\rho_i/\theta)$ is the thermo-chemical potential of the $i$–th species. The parameter $\rho$ describes the total mass density, i.e. $\rho = \sum_{i=1}^{n} \rho_i$, $\kappa(\theta)$ describes the heat conductivity and $r_i$ denotes reaction terms for $i = 1, \ldots, n$. The terms $M_i \nabla (1/\theta)$ and $\sum_{j=1}^{n} M_j \nabla q_j$ describe the Soret and Dufour effect, respectively.

We prescribe the following boundary and initial conditions

$$J_i \cdot \nu = 0, \quad J_e \cdot \nu + \lambda(\theta_0 - \theta) = 0 \quad \text{on } \partial\Omega, \; t > 0, \tag{1.3}$$

$$\rho_i(\cdot, 0) = \rho_i^0, \quad (\rho_i \theta)(\cdot, 0) = \rho_i^0 \theta^0 \quad \text{in } \Omega, \; i = 1, \ldots, n, \tag{1.4}$$

where $\nu$ is the outer normal vector of $\partial\Omega$ and $\theta_0 > 0$, $\lambda \geq 0$ are constant.

We call the diffusion fluxes in equation (1.1) the Fick–Onsager formulation. This formulation comes from the Onsager reciprocal relations [Ons31]. In the isothermal case, it has been shown in [BD23] that the Fick–Onsager form and the usual Maxwell–Stefan formulation which is given by

$$\partial_t \rho_i + \operatorname{div} J_i = r_i, \quad d_i = -\sum_{j=1}^{n} b_{ij} \rho_i \rho_j \left( \frac{J_i}{\rho_i} - \frac{J_j}{\rho_j} \right), \quad i = 1, \ldots, n$$

where $b_{ij} = b_{ji} \geq 0$ for $i, j = 1, \ldots, n$, are equivalent.

As the heat flux, i.e. the term $-\kappa(\theta)\nabla\theta$, is given by Fourier's law, we call the whole system the Maxwell–Stefan–Fourier system in Fick–Onsager form.

Naturally, we assume conversation of the total mass in our system which means, that the sum over the diffusion fluxes $J_i$ as well as the sum over the reaction terms $r_i$ should vanish. Therefore, the total mass density $\rho$ is constant in time. Furthermore, we assume

$$\sum_{i=1}^{n} M_{ij} = 0 \quad \text{for } j = 1, \ldots, n, \text{ and } \sum_{i=1}^{n} M_i = 0. \tag{1.5}$$

### 1.1.2 Mathematical Challenges

We prove two existence results; one for the nondegenerate system and the one for the degenerate case. In the first case, we assume that the diffusion coefficients are symmetric, i.e.

$M_{ij} = M_{ji}$ for all $i, j = 1, \ldots, n$, and that the diffusion matrix $M_{ij}$ is positive semidefinite in the sense, that for $c_M > 0$ holds

$$\sum_{i,j=1}^{n} M_{ij}(\boldsymbol{\rho}, \theta) z_i z_j \geq c_M \left| \Pi \boldsymbol{z} \right|^2 \quad \text{for } \boldsymbol{z} \in \mathbb{R}^n, \, \boldsymbol{\rho} \in \mathbb{R}_+^n, \theta \in \mathbb{R}_+, \tag{1.6}$$

where $\Pi = I - \frac{1}{n} \mathbf{1} \otimes \mathbf{1}$ with $\mathbf{1} = (1, \ldots, 1) \in \mathbb{R}^n$ is the orthogonal projection on $\operatorname{span}\{\mathbf{1}\}^\perp$. Therefore, we do not have coercivity of the diffusion operator, which would be necessary to obtain $H^1$–estimates for the chemical potentials. Furthermore, the equation (1.1) has (without additional assumptions) a singularity in $\rho_i = 0$ and $\theta = 0$. Hence, we have to ensure the positivity of the partial mass densities and the temperature, which is not trivial.

In the second case, we weaken condition (1.6), such that we allow the degeneracy in the partial mass densities. Namely, we assume

$$\sum_{i,j=1}^{n} M_{ij}(\boldsymbol{\rho}, \theta) z_i z_j \geq c_M \sum_{i=1}^{n} \rho_i \left( \Pi \boldsymbol{z} \right)_i^2 \quad \text{for } \boldsymbol{z} \in \mathbb{R}^n, \, \boldsymbol{\rho} \in \mathbb{R}_+^n, \, \theta \in \mathbb{R}_+. \tag{1.7}$$

To overcome these issues, we first use the volume filling assumption $\sum_{i=1}^{n} \rho_i = \rho$ to eliminate one equation. This means we only solve the equation for $\rho_i$, $i = 1, \ldots, n-1$ and obtain $\rho_n$ from the relation $\rho_n = \rho - \sum_{i=1}^{n-1} \rho_i$. This is advantageous, because we obtain positive definiteness of the reduced diffusion matrix $(M_{ij})_{i,j=1}^{n-1}$ under the assumption (1.6). Then, we adapt the techniques of [Jün15] by using the entropy structure of the system. To be more precise, we introduce the mathematical entropy and define the (relative) chemical potential $v_i$ as well as $w$ with $v_i = \log \rho_i - \log \rho_n$ and $w = \log \theta$. By inverting the relation between $v_i$ and $\boldsymbol{\rho}$, we find for $i = 1, \ldots, n-1$

$$\rho_i = \frac{\rho^0 \exp(v_i)}{\sum_{j=1}^{n} \exp(v_j)}.$$

This ensures $0 < \rho_i < \rho^*$ as well as $\theta > 0$. Using an implicit Euler discretization in time, we then solve a regularized problem in variables $v_i$ and $w$. Finally, we derive suitable estimates in form of an entropy inequality and apply the Aubin–Lions compactness lemma to perform the deregularization.

### 1.1.3 State of the Art

We repeat the state of the art which we have given in [HJ21] and extend it by the recent research since the publication.

The isothermal equations were derived from the multi-species Boltzmann equations in the diffusive approximation in [BB21, BGPS13]. The Fick–Onsager form of the Maxwell–Stefan equations was rigorously derived in Sobolev spaces from the multi-species Boltzmann system in [BG20]. The Maxwell–Stefan equations in the Fick–Onsager form, coupled with the momentum balance equation, can be identified as a rigorous second-order Chapman–Enskog approximation of the Euler (–Korteweg) equations for multicomponent fluids; see [HJT19] for the Euler–Korteweg case and [OR20] for the Euler case. The work [BGP19] is concerned with the friction limit in the isothermal Euler equations using the hyperbolic formalism developed by Chen, Levermore, and Liu. A formal Chapman–Enskog expansion

of the stationary non-isothermal model was presented in [TA99]. Another non-isothermal Maxwell–Stefan system was derived in [ABSS20], but the energy flux is different from the expression in (1.2).

The existence analysis of (isothermal) Maxwell–Stefan equations started with the paper [GM98], where the existence of global-in-time weak solutions near the constant equilibrium was proved. A proof of local-in-time classical solutions to Maxwell–Stefan systems was given in [Bot11], and regularity and instantaneous positivity for the Maxwell–Stefan system were shown in [HMPW17]. In [JS13], the entropy or formal gradient-flow structure was revealed, which allowed for the proof of global-in-time weak solutions with general initial data. Maxwell–Stefan systems coupled with the Poisson equation for the electric potential, were analyzed in [JL19].

Alt and Luckhaus [AL83] proved a global existence result for parabolic systems related to the Fick–Onsager formulation. However, their result cannot be applied directly to system (1.1) because of the lack of coerciveness. Moreover, this theory does not yield $L^\infty$ bounds. They are obtained from the technique of [Jün15], but the treatment of Soret and Dufour terms requires some care and is not contained in that work. In [BD23] the relation between Maxwell–Stefan and Fick–Onsager formulation was thoroughly investigated in the isothermal case. All the mentioned results hold if the barycentric velocity vanishes. For non-vanishing fluid velocities, the Maxwell–Stefan equations need to be coupled with the momentum balance. The Maxwell–Stefan equations were coupled with the incompressible Navier–Stokes equations in [CJ15], and the global existence of weak solutions was shown. A similar result can be found in [DD21], where the incompressibility condition was replaced by an artificial time derivative of the pressure and the limit of vanishing approximation parameters was performed. Coupled Maxwell–Stefan and compressible Navier–Stokes equations were analyzed in [BD21], and the local-in-time existence analysis was performed. A global existence analysis for a general isothermal Maxwell–Stefan–Navier–Stokes system was performed in [DDGG20]. For the existence analysis of coupled stationary Maxwell–Stefan and compressible Navier–Stokes–Fourier systems, we refer to [BJPZ22, GPZ15, PP17]. In [BJPZ22], temperature gradients were included in the partial mass fluxes, but only the stationary model was investigated. In [Dru22], a Navier–Stokes–Fick–Onsager–Fourier system is discussed, where the results from [BD21] are generalized to the nonisothermal case. The paper [FHKM22] discusses a more general class of nonisothermal reaction–diffusion systems including the Soret and Dufour effect and proves the global existence of renormalised solutions. Furthermore, the most recent paper [JG23] is based on the model, which we discuss in [HJ21], i.e. this thesis, and improves on the modeling regarding the thermodynamics. However, the global existence of solutions can still be proved with the same techniques and (partly) the same estimates as in this thesis.

## 1.2 About Biofilms

Biofilms are accumulations of microorganisms which grow on surfaces and produce extracellular polymeric substance (EPS) ( [PMCW11]). The EPS can be understood as a layer of slime and varies wildly depending on the underlying microorganisms, see for in-

stance [PMCW11] for a description of the composition for gram–negative bacteria. The production of EPS enables a protected mode of growth for the microorganisms, which means the Biofilm structure is more resistant against antimicrobials [Poz18].

The biofilm development can be divided into several stages (see [KD10]). At first, free living (planktonic) microorganisms attach to a surface and become sessile. These cells then connect by producing EPS and thus forming a growing biofilm. Lastly, cells detach from the biofilm due to different effects; for instance, erosion, mechanical stress ( [KD10]) or quorum sensing. Quorum sensing describes the communication between cells via signal molecules [EHKE15]. These signal–molecules, which are also called autoinducers, are able to cause a detachment of cells in the biofilm. This can lead to a dispersion of microorganisms, which can lead to the development of new biofilm colonies, see for instance [SEL14]. In this work, we do not consider erosion or mechanical stress, but we consider in section 1.3 rather a model, which accounts quorum sensing as cause for a detachment of cells.

As biofilms can be prevelant on almost any surface in a moist environment, they are of major importance in nature, food industry and medicine, see for instance [FAG09, HSCS04, SN16, Bry08].

Therefore, we want to emphasize that infections with biofilms, as for instance on catheters by Staphyloccucus aureus bacteria, are highly problematic, since the biofilm structure gives an improved protection of antibacterial treatment [HSCS04].

## 1.3 A Quorum Sensing induced Biofilm Model

We discuss a system of nonlinear partial differential equations which models the biofilm growth including quorum sensing effects. It is an extension to the biofilm growth model suggested in [EPL01] and describes the growth of a biofilm dependent on the nutrient, signal molecules and the cells which got detached due to the quorum sensing effect of the signal molecules. The model was first suggested in [EHKE15] and then mathematically analyzed in [ESE17]. The aim in this thesis is, to define an implicit Euler finite–volume scheme for the model analyzed in [ESE17] and to prove the existence of discrete solutions which preserve the $L^\infty$–bounds of the model, as well as the convergence towards a weak solution.

### 1.3.1 The Model Equations

We begin by introducing the model of [ESE17]. The biofilm is modeled by the biomass fraction $M(x,t)$. We say biomass, as $M(x,t)$ describes technically the EPS, which includes the microorganisms as well as other substances. The nutrient concentration is modeled by $S(x,t)$. The parameter $A(x,t)$ describes the autoinducer, i.e. the signal molecule, which induces the quorum sensing effect. Lastly, the dispersed cells are modeled by $N(x,t)$. Then, $M, N, S, A$ satisfy the scaled diffusion equations

$$\partial_t M - d_1 \operatorname{div}(f(M)\nabla M) = g_1(M,S,A) \quad \text{in } \Omega, \ t > 0, \tag{1.8}$$

$$\partial_t N - d_2 \Delta N = g_2(M,N,S,A), \quad \text{in } \Omega, \ t > 0, \tag{1.9}$$

$$\partial_t S - d_3 \Delta S = g_3(M,N,S), \quad \text{in } \Omega, \ t > 0, \tag{1.10}$$

$$\partial_t A - d_4 \Delta A = g_4(M, N, A), \quad \text{in } \Omega, \ t > 0, \tag{1.11}$$

and the initial and boundary conditions

$$M(0) = M^0, \ N(0) = N^0, \ S(0) = S^0, \ A(0) = A^0 \quad \text{in } \Omega, \tag{1.12}$$

$$M = M^D, \ N = 0, \ S = 1, \ A = 0 \quad \text{on } \partial\Omega, \ t > 0, \tag{1.13}$$

where $\Omega \subset \mathbb{R}^d$ $(d \geq 1)$ is a bounded domain and $M^D \in (0, 1)$ is a constant.

**Remark 1.** *In the literature, mostly homogeneous Dirichlet boundary conditions have been used but for the numerical analysis, we need nonhomogeneous boundary conditions since the introduction of the entropy variable requires $M^D$ to be nonzero. However, we may $M^D = \gamma$ and pass the limit $\gamma \to 0$ to cover homogeneous conditions as well.*

The source terms used in [ESE17] describe the nutrient consumption, the dead of biomass and dispersed cells, the production of signal molecules and detachment of biomass.

$$g_1(M, S, A) = \frac{S}{k_1 + S} M - k_2 M - \eta \left( \frac{A^n}{1 + A^n} \right) M, \tag{1.14}$$

$$g_2(M, N, S, A) = \frac{S}{k_1 + S} N - k_2 N + \eta \left( \frac{A^n}{1 + A^n} \right) M, \tag{1.15}$$

$$g_3(M, N, S) = -\frac{\mu S}{k_1 + S} (M + N), \tag{1.16}$$

$$g_4(M, N, A) = -\lambda A + \left[ \alpha + \beta \frac{A^n}{1 + A^n} \right] (M + N). \tag{1.17}$$

where $k_1, k_2, \alpha, \beta, \eta, \mu > 0$ and $n > 1$.

The growth of biomass and dispersed cells in (1.14) and (1.15) respectively is controlled by the nutrient availability and described by the monod kinetic growth term $S/(k_1 + S)$. Furthermore, the biomass and the dispersed cells both die with rate $k_2$. However, the last term of the source terms $g_1$ and $g_2$ describes the detachment through quorum sensing: cells from the biomass get detached with rate $\eta A^n/(1 + A^n)$ and become dispersed cells. The nutrient consumption by the biomass $M$ and dispersed cells $N$ is described in the source term $g_3$. The signal molecule/autoinducer has a decay rate of $\lambda$ and is produced with rate $\alpha + \beta A^n/(1 + A^n)$.

**Remark 2.**   *(i) The original model in [EHKE15] also contains the effect of re–attachment of dispersed cells to the biomass. According to [ESE17], the effect of re–attachment is negligible, due to why it is not treated in the analysis.*

*(ii) Furthermore, as pointed out in [ESE17] as well, without a signal molecule production/activity, i.e. $\alpha = \beta = \eta = 0$ we recover the originally suggested model from [EPL01]. Thus, the model (1.8)–(1.17) be considered as a generalization.*

The diffusion term $f(M)$ is chosen, as suggested in [EPL01] and used in [ESE17], as

$$f(M) = \frac{M^b}{(1 - M)^a}, \quad \text{where } a > 1, \ b > 0. \tag{1.18}$$

We can rewrite the diffusion operator in equation (1.8) as $\mathrm{div}(f(M)\nabla M) = \Delta F(M)$, where

$$F(M) = \int_0^M f(s)ds, \quad M \geq 0. \tag{1.19}$$

Heuristically, the superdiffusion singularity prevents the biomass fraction to exceed its maximum value 1 as the diffusion gets larger and "spreads" the biomass, while the porous–medium degeneracy leads to finite speed of propagation. For more details regarding the modeling, we refer to [EPL01, EHKE15, ESE17].

### 1.3.2 Mathematical Challenges

The main difficulty of the analysis is the degenerate-singular diffusion term. On the continuous level, the authors of [ESE17] proved that the choice of the initial value $M^0$ such that $\|M^0\|_{L^\infty(\Omega)} < 1 - \rho$ for some $\rho \in (0,1)$ guarantees the existence of $\delta > 0$ such that $M \leq 1 - \delta$ almost everywhere in $\Omega$ ( [ESE17, Lemma 3.3]). This bound is proved by the use of a comparison principle, which we could not adapt to the discrete case. We overcome this challenge by introducing an entropy variable $W$ of the form

$$W^\varepsilon = F(M) - F(M^D) + \varepsilon \log \frac{M^\varepsilon}{M^D} \tag{1.20}$$

and regularize the (discretized) system by adding higher order terms of the form $\varepsilon(W^\varepsilon - \Delta W^\varepsilon)$, where $\varepsilon$ is the parameter for the regularization. Solving the regularized problem in $W^\varepsilon$ and using the invertibility of $(0,1) \to \mathbb{R}$, $M^\varepsilon \mapsto W^\varepsilon$, we find $0 < M^\varepsilon < 1$.

### 1.3.3 Finite–Volume Methods

We give a short introduction to finite–volume methods. To this end, we focus solely on the discretization of the biomass equation (1.8), as it is clear that the discretization for the other equations works analogously. We want to mention on a lighter note, that finite–volume methods is a broad research topic and we just scratch the surface with our simplified motivation. Thus, we refer to the finite–volume monument [EGH00] for a detailed introduction into finite–volume methods. To discretize the scheme, we first discretize with an implicit euler method in time. Therefore, we replace the continuous time derivative by

$$\partial_t M \approx \frac{M^k - M^{k-1}}{\Delta t}$$

for $k = 1, \ldots, N_T$, where $N_T$ denotes the number of time steps and $\Delta t$ the size of the time step. Then, the idea of the finite–volume method is, to partition the domain $\Omega$ such that $\bigcup_{K \in \mathcal{T}} K = \Omega$, where $\mathcal{T}$ denotes the set of control volumes and $K$ denotes the control volumes/cells. We integrate equation (1.8) over one control volume $K$ and formally apply the divergence theorem to obtain

$$\int_K \frac{M^k - M^{k-1}}{\Delta t}dx - \sum_{\sigma \in \mathcal{E}_K} \int_\sigma d_1 \nabla F(M^k) \cdot \nu_{K,\sigma}ds = \int_K g_1(M^k, S^k, A^k)dx, \tag{1.21}$$

where $\sigma$ denotes an edge, $\mathcal{E}_K$ denotes the set of edges of $K$ and $\nu_{K,\sigma}$ denotes the outside normal vector of control volume/cell $K$ on edge $\sigma$. The goal is, to rewrite the equation in terms of $M_K^k$, where $M_K^k = \mathrm{m}(K)^{-1} \int_K M^k(x) dx$, where $\mathrm{m}(K)$ denotes the volume of the cell $K$. For the first term of equation (1.21), this is rather obvious. Furthermore, we replace $(M^k, S^k, A^k)$ in the source term by the averages over the cell $K$ to obtain

$$\int_K g_1(M^k, S^k, A^k) dx \approx \mathrm{m}(K) g_1(M_K^k, S_K^k, A_K^k).$$

It remains to approximate the integral over the boundary of cell $K$. To this end, we denote with $\mathcal{E}_K$ the edges of cell $K$ and distinguish two cases:

(i) The edge $\sigma \in \mathcal{E}_K$ is separating the cells $K$ and $L \in \mathcal{T}$, which we denote as $\sigma = K \mid L$. In this case, we assume that the edge $\sigma$ is orthogonal to the straight line which connects the middle points $x_K$ and $x_L$ of cells $K$ and $L$, respectively.

$$- \int_\sigma d_1 \nabla F(M^k) \cdot \nu_{K,\sigma} ds \approx -d_1 \frac{\mathrm{m}(\sigma)}{\mathrm{d}_\sigma} \left( F(M_L^k) - F(M_K^k) \right) =: \mathcal{F}_{M,K,\sigma}^k, \qquad (1.22)$$

where $\mathrm{m}(\sigma)$ denotes the length of edge $\sigma$ and $\mathrm{d}_\sigma$ the distance between $x_K$ and $x_L$.

(ii) The edge $\sigma \in \mathcal{E}_K$ is an exterior edge, i.e. $\sigma \subset \partial\Omega$. In this case we approximate the integral over $\sigma$ by

$$- \int_\sigma d_1 \nabla F(M^k) \cdot \nu_{K,\sigma} ds \approx -d_1 \frac{\mathrm{m}(\sigma)}{\mathrm{d}_\sigma} \left( F(M^D) - F(M_K^k) \right) =: \mathcal{F}_{M,K,\sigma}^k, \qquad (1.23)$$

where $\mathrm{d}_\sigma$ denotes the distance between the middle point $x_K$ and the edge $\sigma$.

The approximation in equations (1.22)–(1.23) is called Two–Point flux approximation. We want to point out, that the orthogonality assumption for the edges $\sigma$ which we made in case 1.3.3 is crucial for the consistency of the approximation, i.e. that the truncation error on the flux is of the samer order as the maximum length of the edges of the mesh (see [EGH00, Example 1.2]). Summarized, we reach the following finite–volume approximation for equation (1.8):

$$\frac{\mathrm{m}(K)}{\Delta t} \left( M_K^k - M_K^{k-1} \right) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{M,K,\sigma}^k = \mathrm{m}(K) g_1(M_K^k, S_K^k, A_K^k).$$

### 1.3.4 State of the Art

Due to the importance of biofilm analysis, there are many different mathematical models for biofilm growth. We do not attempt to list all of them but rather give a rough overview. The mathematical modeling of biofilms started in the 1980s with the work of [RM80] which focused on the biofilm dynamics at steady states. In [WG86], a one dimensional model including a transport equation for biomass and describing the evolution of the biofilm thickness was suggested (see [WZ10] for further details and references). The work [PvLH99] suggested a multidimensional model which implements the existence

of a sharp front between the biomass and the fluid. In this model, the domain is separated into a liquid and a solid area. The solid area is characterized by the positivity of biomass density while the liquid area is characterized by the absence of biomass. The flow velocity in the liquid area satisfies the incompressible Navier–Stokes equation. The growth of the biomass is controlled by a reaction term. The equations for this model do not contain spreading of the biomass. However, the effect of spreading itself is considered as "redistribution in space according to discrete rules" ( [PvLH99]). For more details, we refer to [PvLH99, PvLH98b, PVLH98a]. In [EPL01], the authors extend [PvLH99] by suggesting a new equation to model the biomass growth, which also considers the spatial distribution by including a diffusion flux. For this model, the existence and uniqueness of global weak solutions has been shown in [EZE09] for the hydrostatic case, i.e. without coupling with the incompressible Navier–Stokes equation. The (hydrostatic) model has been extended in several ways, for instance by adding a nutrient taxis term to the biomass equation in [EEWZ14] or considering quorum sensing effects in [EHKE15]. The global existence of weak solutions has been shown in [ESE17]. Another approach is the modeling of mixing effects by considering multiple biofilm species, as for instance in [RSE15]. For this model, global existence of weak solutions has been shown in [DMZ19], while a finite–volume scheme was developed in [DJZ21] including numerical analysis. However, this model does not consider an equation for the nutrient. A finite–volume method of the model of [EPL01] was considered in [AES18], but without containing numerical analysis. Other variants of the above biofilm model are possible. For instance, the authors in [HES22] considered a PDE–ODE system which contains the equation for biomass growth from [EPL01] and an ODE for the nutrient consumption. After a spatial discretization, random cell attachment is considered and numerically simulated for the resulting ODE.

A completely different approach to the modeling of biofilm growth was chosen in [FHX14], where the authors models the biofilm evolution as a two–phase free boundary problem. In this model, the fluid outside of the biomass is considered as an incompressible viscous fluid, while within the biofilm, a mixture of two fluids is assumed. The interface between both areas is representing the free boundary. Lastly, we mention the model of [WZ12]. In this model, the biofilm is modeled as a fluid–mixture consisting of the biomass and solvent containing dissolved nutrient substrate. The equation for the fluid is coupled with a reaction–diffusion equation, which describes the substrate concentration. The evolution of the biofilm is then governed by a chemical potential coming from the extended Flory–Huggins energy (see [ZCW08a, ZCW08b, WZ12]), leading to a Cahn–Hilliard Type equation.

## 1.4 A Cahn–Hilliard Type system modeling Biofilm Growth

We discuss a biofilm growth model motivated by [WZ12]. As mentioned in section 1.3.4, the biofilm is modeled as a fluid mixture, which leads to a Cahn–Hilliard type equation for the biomass growth. The aim is to establish an existence result for this kind of biofilm growth model.

The Cahn–Hilliard equation was developed to describe a phase separation between two components [NC08]. In our case, the two components are the biomass fraction $u$ and the solvent fraction $u_s$. Before going into more details, we introduce the model equations.

### 1.4.1 The Model Equations

The (modified) version of the model in [WZ12] is given by

$$\partial_t v - \operatorname{div}((1-u)\nabla v) = g(u,v), \tag{1.24}$$

$$\partial_t u - \operatorname{div}(M(u)\nabla\mu) = h(u,v), \tag{1.25}$$

$$\mu = -\Delta u + f'(u) \quad \text{in } \Omega, \ t > 0, \tag{1.26}$$

where $\Omega \subset \mathbb{R}^d$ ($d \geq 1$) is a bounded domain. We impose no–flux boundary conditions

$$(1-u)\nabla u \cdot \nu = M(u)\nabla\mu \cdot \nu = \nabla u \cdot \nu = 0 \quad \partial\Omega, t > 0. \tag{1.27}$$

The initial conditions are

$$u(0) = u^0, \quad v(0) = v^0. \tag{1.28}$$

The mobility is chosen as

$$M(u) = u(1-u), \tag{1.29}$$

although more general choices are possible (see Remark 34). The function $f$ describes the Flory–Huggins mixing free energy

$$f(u) = \frac{1}{N}u\log u + (1-u)\log(1-u) + \lambda u(1-u), \tag{1.30}$$

$$f'(u) = \frac{1}{N} + \lambda + \frac{1}{N}\log u - \log(1-u) - 1 - 2\lambda u \tag{1.31}$$

where $N > 0$ is the generalized polymerization index and $\lambda > 0$ describes the Flory–Huggins mixing parameter. We define the reaction terms as

$$g(u,v) = -ug_0(v), \ h(u,v) = u(1-u)h_0(v), \tag{1.32}$$

where $g_0$ and $h_0$ are non decreasing, continuous functions.

Equation (1.25) describes the evolution of the biomass fraction as per a degenerated Cahn–Hilliard equation with nonlinear source term. The degeneracy comes from the degenerate mobility $M(0) = M(1) = 0$. Furthermore, the equation is singular due to the derivative of the Flory–Huggins mixing free energy (1.31). As mentioned in the introduction, the Cahn–Hilliard equation describes the phase separation between two components. The second component is the solvent fraction, which is implicitly given by the volume filling assumption $u + u_s = 1$, such that $u_s = 1 - u$. The case $u = 0$ is corresponding from model point of view to the abscence of biomass, i.e. only solvent is present. The case $u = 1$ on the other hand would imply the abscence of solvent and as such the presence of purely biomass. In [ZCW08a], this case is excluded as a modeling assumption, as it would imply the existence of dry biomass. In the modified equations, we allow the "dry–biomass" case in the mathematical analysis, as we can not guarantee $u < 1$ with the chosen mobility. However, an appropriate choice of the mobility excludes the case $u = 1$ (see Remark 34). Equation (1.24) describes the substrate concentration, which is necessary for the biomass

growth. The substrate consumption is only possible in the solvent (as it needs to be dissolved) and then consumed with rate $h_0(v)$. The factor $u(1-u)$ in $h$ ensures, that there is no growth in the pure solvent/pure biomass case. This implies, that in the pure solvent case bacterial mass can not come out of nowhere, while in the pure biomass case the substrate can not be dissolved in the solvent to cause further growth. The function $g$ describes the consumption of the substrate and is chosen as in [WZ12]. From a modeling point of view, this is slightly inaccurate in the case $u = 1$: The growth stops, the solvent is not present and the substrate can not be dissolved in solvent. Thus, the consumption of the substrate should also vanish. This can be easily fixed by also adding a factor $1-u$ in the definition of $g$ or changing the mobility to ensure $u < 1$, see Remark 34. However, as we want to compare our results numerically to the model of [WZ12], we decided to keep the function $g$ as defined in [WZ12], which also remains flexible in the choice whether to adjust the mobility or the source term.

Since our model is motivated by [WZ12] but yet different, we briefly remark the differences between the two models:

(i) We neglected the velocities of biomass and solvent,

(ii) We added the solvent fraction $u_s = 1 - u$ in the production rate $h$ and the mobility $M$,

(iii) We neglect the elastic energy to mathematically simplify the definition of the chemical potential $\mu$,

(iv) We neglect the factor of the solvent fraction in the time derivative of equation (1.24).

Our analysis works as well with given velocities with bounded divergence. The addition of the solvent fraction in the mobility and production rate seems however mathematically necessary to achieve an existence result. To be more precise, we need a cancellation in $M(u)f''(u)$ to identify the weak limit and does not seem possible without the additional factor. The negligence of the solvent fraction in the time derivative removes a degeneracy which we were not able to treat.

### 1.4.2 Mathematical Challenges

As there is no general maximum principle for equations of fourth order, it is difficult to prove the $L^\infty$–bounds $0 \le u \le 1$ for the biomass fraction. Furthermore, due to the degeneracies in the mobilities of equations (1.24) and (1.25), we do not obtain an $H^1$ estimate for the chemical potential and the substrate fraction, respectively.

As a consequence, we can not expect strong convergence for approximate solutions of the substrate equation. We overcome these issues in the following way: First, we truncate the mobility $M(u)$, the Flory–Huggins mixing energy $f(u)$ and the source terms. Furthermore, we add a higher order term of the form $\kappa\Delta v$ for the substrate equation. Following then [EG96], we use a Galerkin method and obtain uniform estimates by entropy/energy inequalities. Using the Stampacchia method, we prove the bounds for the substrate equation, i.e. $0 \le v \le 1$. With the help of an entropy inequality and the Lemma of Fatou we conclude the $L^\infty$–bound for the biomass fraction, i.e. $0 \le u \le 1$ in the deregularization

limit. However, there are still some difficulties: Since we do not have $H^1$ estimates for the chemical potential and the substrate concentration, we can only identify the limit in the sense of $L^2(0, T; H^1(\Omega)')$. Furthermore, due to lack of strong convergence for the substrate in the deregularization, it proves difficult to identify the limit in the nonlinear source terms. We overcome this by using a Minty–Browder argument.

### 1.4.3 State of the Art

Since we gave already a state of the art regarding biofilm models (see section 1.3.4), we focus on the analytical aspects of Cahn–Hilliard equations in this section. To this end, we repeat the relevant part of the state of the art which we have given in [HJ23] and extend it by references which caught our attention only after submission as well as other interesting results. The first existence analysis of Cahn–Hilliard equations was given in [Yin92] in one space dimension and in [EG96] in several space dimensions. Most of the analytical results on the Cahn–Hilliard equations do not contain reaction terms. Moreover, if reaction terms are included in the Cahn–Hilliard model, mostly nondegenerate mobilities are chosen; see, e.g., [AKK11, CMZ14, GL16]. When the gradient term in the free energy is replaced by a nonlocal spatial interaction energy, degenerate mobilities (and singular potentials) can be treated [Fri21, IM18]. As per our knowledge, there are few papers which consider degenerate mobilities combined with source terms, for instance [FLR17, MR15]. However, the setting is different than in the model which we consider, as both papers examine a nonlocal variant the Cahn–Hilliard equation. The only work which considers degenerate mobility in combination with source terms is [Ebe20]: However, our work is different in two important details: In [Ebe20, Chapter 7.3], upper bounds can not be proved and in [Ebe20, Chapter 7.5], the source terms contain the chemical potential. Furthermore, we have a second degeneracy in the coupled reaction–diffusion equation which causes additional difficulties. Some results describe a connection from nonlocal to local models. For instance, in [CES23] the convergence to local solutions on the torus without source terms is proved. In [MRST19] the convergence of solutions to a nonlocal Cahn–Hilliard equation with constant mobility and without source terms to solutions of the local Cahn–Hilliard equation is proved. The connection between nonlocal models and local models considering degeneracies and source terms seems yet unclear.

## 1.5 Outline of the Thesis

In this section, we give an overview about the structure of the thesis.

Chapter 2 is considering the Maxwell–Stefan–Fourier system in Fick–Onsager form.

- We explain the modeling of the system in section 2.3.

- We prove that the Maxwell–Stefan formulation leads to the Fick–Onsager form for a specific choice of coefficients (Proposition 5).

- We formulate the existence theorems for the nondegenerate case (Theorem 3) and the degenerate case (Theorem 4).

- We give the proof for theorem 3 in section 2.4 and the proof for theorem 4 in section 2.5 respectively.

The results in this chapter are based on the research collaboration with Ansgar Jüngel (TU Wien) and has been under the title *Analysis of Maxwell–Stefan systems for heat conducting fluid mixtures* ( [HJ21]).

In Chapter 3, we present the details about the finite–volume scheme/numerical analysis for the quorum sensing biofilm model.

- In Section 3.1, we formulate the finite–volume scheme.

- We formulate the existence for discrete solutions (Theorem 13), the uniqueness of discrete solutions (Theorem 14) and the convergence of discrete solutions towards weak solutions of the (continuous) equations (Theorem 15).

- We present the proofs for the existence, uniqueness and convergence in sections 3.2, 3.3 and 3.5, respectively.

- We present some numerical experiments in section 3.6.

The results of this chapter are based on a research collaboration with Ansgar Jüngel (TU Wien) and Antoine Zurek (UTC) published under *Analysis of a finite–volume scheme for a single–species biofilm model* ( [HJZ23]).

In chapter 4, we provide the existence proof for the Cahn–Hilliard Type biofilm model.

- We formulate the existence theorem for global weak solutions (Theorem 22) and explain the key ideas.

- In sections 4.2–4.4 we provide the details for the existence proof.

- In section 4.5, we discretize the system by a BDF2 discretization in time and a finite–volume discretization in space.

- We present some numerical experiments and compare our results to the model of [WZ12].

The results of this chapter are based on a research collaboration with Ansgar Jüngel (TU Wien) and are submitted for publication under the title *Existence Analysis for a reaction–diffusion Cahn–Hilliard–Type system with degenerate mobility and singular potential modeling biofilm growth* ( [HJ23]).

Lastly, to conclude the thesis, we give a short discussion of our results and give an outlook to further possible research directions in chapter 5.

# 2 Analysis for the Maxwell–Stefan–Fourier system

The results in this chapter have been published in [HJ21].

In this chapter, we provide the mathematical details for the Maxwell–Stefan Fourier system in Fick–Onsager form. To this end, we first present the main results, namely the existence theorem in the nondegenerate and the degenerate case in section 2.1. We present the mathematical ideas used in the existence proof in section 2.2. In section 2.3, we explain the derivation of the Maxwell–Stefan Fourier system in Fick–Onsager form and prove, that the Maxwell–Stefan formulation implies the Fick–Onsager form for a suitable choice of $M_{ij}$, $M_i$ and $d_i$. In section 2.4 we then give the full proof for the existence in the nondegenerate case, and in section 2.5 we describe, which steps in the proof of section 2.4 changes to obtain the degenerate existence result.

## 2.1 Main Results

Before stating our main results, we impose the following assumptions: We impose the following assumptions:

(H1)  *Domain:* $\Omega \subset \mathbb{R}^3$ is a bounded domain with a Lipschitz continuous boundary.

(H2)  *Data:* $\theta^0 \in L^\infty(\Omega)$, $\inf_\Omega \theta^0 > 0$, $\theta_0 > 0$, $\lambda \geq 0$; $\rho_i^0 \in H^1(\Omega) \cap L^\infty(\Omega)$ satisfies $0 < \rho_* \leq \rho_i^0 \leq \rho^*$ in $\Omega$ for some $\rho_*, \rho^* > 0$.

(H3)  *Diffusion coefficients:* For $i, j = 1, \ldots, n$, the coefficients $M_{ij}, M_j \in C^0(\mathbb{R}_+^n \times \mathbb{R}_+)$ satisfy (1.5) and $M_{ij}$, $M_i/\theta$ are bounded functions.

(H4)  *Heat conductivity:* $\kappa \in C^0(\mathbb{R}_+)$ and there exist $c_\kappa, C_\kappa > 0$ such that for all $\theta \geq 0$,
$$c_\kappa(1 + \theta^2) \leq \kappa(\theta) \leq C_\kappa(1 + \theta^2).$$

(H5)  *Reaction rates:* $r_1, \ldots, r_n \in C^0(\mathbb{R}^n \times \mathbb{R}_+) \cap L^\infty(\mathbb{R}^n \times \mathbb{R}_+)$ satisfies $\sum_{i=1}^n r_i = 0$ and there exists $c_r > 0$ such that for all $\boldsymbol{q} \in \mathbb{R}^n$ and $\theta > 0$,
$$\sum_{i=1}^n r_i(\Pi\boldsymbol{q}, \theta)q_i \leq -c_r|\Pi\boldsymbol{q}|^2.$$

The bounds on $\rho^0$ in Hypothesis (H2) are needed to derive the positivity and boundedness of the partial mass densities. In the example presented in Section 2.3, the coefficients $M_{ij}$

and $M_i/\theta$ depend on $\rho_i$; since we prove the existence of $L^\infty$ solutions $\rho_i$, the functions $M_{ij}$ and $M_i$ are indeed bounded, as required in Hypothesis (H3). The growth condition for the heat conductivity in Hypothesis (H4) is used to derive higher integrability of the temperature, see (2.9), which allows us to treat the heat flux term. If $\lambda = 0$, we can impose the weaker condition $\kappa(\theta) \geq c_\kappa \theta^2$. Hypothesis (H5) is satisfied for the reaction terms used in [DDGG20]. The bound for $\sum_{i=1}^n r_i q_i$ gives a control on the $L^2(\Omega)$ norm of $\Pi\boldsymbol{q}$. Together with the estimates for $\nabla(\Pi\boldsymbol{q})$ from (2.7), we are able to infer an $H^1(\Omega)$ estimate for $\Pi\boldsymbol{q}$. A more natural $L^2(\Omega)$ bound for $\boldsymbol{q}$ may be derived under the assumption that the total initial density does not lie on a critical manifold associated to the reaction rates; we refer to [DDGG20, Theorem 11.3] for details. Vanishing reaction rates are allowed in Theorem 4 below.

**Theorem 3** (Existence). *Let Hypotheses (H1)–(H5) hold, let $(M_{ij})$ satisfy (1.6), and let $T > 0$. Then there exists a weak solution $(\boldsymbol{\rho}, \theta)$ to (1.1)–(1.4) satisfying $\rho_i > 0$, $\theta > 0$ a.e. in $\Omega_T$,*

$$\rho_i \in L^\infty(\Omega_T) \cap L^2(0, T; H^1(\Omega)) \cap H^1(0, T; H^2(\Omega)'), \tag{2.1}$$

$$v_i \in L^2(0, T; H^1(\Omega)), \quad (\Pi\boldsymbol{q})_i \in L^2(0, T; H^1(\Omega)), \tag{2.2}$$

$$\theta \in L^2(0, T; H^1(\Omega)) \cap W^{1,16/15}(0, T; W^{2,16}(\Omega)'), \quad \log\theta \in L^2(0, T; H^1(\Omega)); \tag{2.3}$$

*where $v_i = \log(\rho_i/\rho_n)$ and $(\Pi\boldsymbol{q})_i = v_i - \sum_{j=1}^n v_j/n$ for $i = 1, \ldots, n$; it holds that*

$$\int_0^T \langle \partial_t \rho_i, \phi_i \rangle dt + \int_0^T \int_\Omega \left( \sum_{j=1}^{n-1} M_{ij} \nabla v_j - \frac{M_i}{\theta} \nabla \log\theta \right) \cdot \nabla \phi_i \, dx dt = \int_0^T \int_\Omega r_i \phi_i \, dx dt, \tag{2.4}$$

$$\int_0^T \langle \partial_t(\rho\theta), \phi_0 \rangle dt + \int_0^T \int_\Omega \kappa(\theta) \nabla\theta \cdot \nabla\phi_0 \, dx dt + \int_0^T \int_\Omega \sum_{j=1}^{n-1} M_j \nabla v_j \cdot \nabla\phi_0 \, dx dt \tag{2.5}$$

$$= \lambda \int_0^T \int_{\partial\Omega} (\theta_0 - \theta)\phi_0 \, dx ds$$

*for all $\phi_1, \ldots, \phi_n \in L^2(0, T; H^1(\Omega))$, $\phi_0 \in L^\infty(0, T; W^{1,\infty}(\Omega))$ with $\nabla\phi_0 \cdot \nu = 0$ on $\partial\Omega$, and $i = 1, \ldots, n$; and the initial conditions (1.4) are satisfied in the sense of $H^2(\Omega)'$ and $W^{2,16}(\Omega)'$, respectively.*

The weak formulation can be written in various variable sets since

$$\sum_{j=1}^{n-1} M_{ij} \nabla v_j = \sum_{j=1}^n M_{ij} \nabla(\Pi\boldsymbol{q})_j = \sum_{j=1}^n M_{ij} \nabla q_j,$$

$$\sum_{j=1}^{n-1} M_j \nabla v_j = \sum_{j=1}^n M_j \nabla(\Pi\boldsymbol{q})_j = \sum_{j=1}^n M_j \nabla q_j,$$

whenever the corresponding variables are defined. Thus, our definition of a weak solution is compatible with (1.1)–(1.2). The proof is based on a suitable approximate scheme, uniform bounds coming from entropy estimates, and $H^1(\Omega)$ estimates for the partial mass densities.

More precisely, we use two levels of approximations. First, we replace the time derivative by an implicit Euler discretization to overcome issues with the time regularity. Second, we add higher-order regularizations for the thermo-chemical potentials and the logarithm of the temperature $w = \log \theta$ to achieve $H^2(\Omega)$ regularity for these variables. Since we are working in three space dimensions, we conclude $L^\infty(\Omega)$ solutions, which are needed to define properly $\rho_i = \exp(w + q_i)$.

A priori estimates are deduced from a discrete version of the entropy inequality (2.7). They are derived from the weak formulation by using $v_i$ and $e^{-w_0} - e^{-w}$ as test functions, where $w_0 = \log \theta_0$. The entropy structure is only preserved if we add additionally a $W^{1,4}(\Omega)$ regularization and some lower-order regularization in $w$. The properties for the heat conductivity allow us to obtain estimates for $\theta$ in $H^1(\Omega)$ and for $\nabla \log \theta$ in $L^2(\Omega)$. Property (1.6) provides gradient estimates for $\boldsymbol{v}$ and, in view of (2.8), also for $\boldsymbol{\rho}$.

Condition (1.6) provides a control on the relative thermo-chemical potentials $v_i$, but it excludes the dilute limit, i.e. situations when the mass densities vanish. This situation is included in [Dru21], which deals with the isothermal case. We are able to replace condition (1.6) by a degenerate one, which allows for dilute mixtures:

$$\sum_{i,j=1}^n M_{ij}(\boldsymbol{\rho}, \theta) z_i z_j \geq c_M \sum_{i=1}^n \rho_i (\Pi \boldsymbol{z})_i^2 \quad \text{for } \boldsymbol{z} \in \mathbb{R}^n, \; \boldsymbol{\rho} \in \mathbb{R}_+^n, \; \theta \in \mathbb{R}_+. \tag{2.6}$$

This corresponds to "degenerate" diffusion coefficients $M_{ij}$; see Section 2.3 for a motivation. Although this hypothesis seems to complicate the problem, there are two advantages. First, it allows us to derive a gradient bound for $\rho_i^{1/2}$, and second, it helps us to avoid the bound from $r_i$ in Hypothesis (H5). In fact, we may assume that $r_i = 0$.

**Theorem 4** (Existence, "degenerate" case)**.** *Let condition* (2.6) *be satisfied. Moreover, let Hypotheses (H1)–(H4) hold for $T > 0$ and additionally, $(\rho_i^0)^{1/2} \in H^1(\Omega) \cap L^\infty(\Omega)$, $M_{ij}/\rho_j$ and $M_j/\rho_j$ are bounded, $r_i = 0$ for all $i, j = 1, \ldots, n$. Then there exists a weak solution $(\boldsymbol{\rho}, \theta)$ to* (1.1)–(1.4) *satisfying $\rho_i \geq 0$, $\theta > 0$ a.e. in $\Omega_T$,* (2.1)*,* (2.3)*, and the weak formulation* (2.4)–(2.5) *with, respectively,*

$$\sum_{i=1}^n \frac{M_{ij}}{\rho_j} \nabla \rho_j, \quad \sum_{i=1}^n \frac{M_i}{\rho_i} \nabla \rho_i \quad \text{instead of} \quad \sum_{i=1}^{n-1} M_{ij} \nabla v_j, \quad \sum_{i=1}^{n-1} M_i \nabla v_i.$$

## 2.2 Mathematical Ideas

In this section, we describe the key ideas of the proofs to theorem 3 and 4. As mentioned in the introduction, we want to use the entropy structure of the system.

We describe the main ideas for the first existence result. We use the mathematical entropy

$$h = \sum_{i=1}^n \rho_i (\log \rho_i - 1) - \rho \log \theta.$$

Introducing the relative thermo-chemical potentials $v_i = \partial h / \partial \rho_i - \partial h / \partial \rho_n = q_i - q_n$ for $i = 1, \ldots, n$ and interpreting $h$ as a function of $(\boldsymbol{\rho}', \theta)$, a formal computation (which is

made precise for an approximate scheme; see (2.20)) shows that

$$\frac{d}{dt}\int_\Omega h(\boldsymbol{\rho}',\theta)dx + \frac{c_M}{2}\int_\Omega\left(\frac{1}{n}|\nabla\boldsymbol{v}|^2 + |\nabla\Pi\boldsymbol{q}|^2\right)dx$$

$$+ \int_\Omega \kappa(\theta)|\nabla\log\theta|^2 dx + \lambda\int_{\partial\Omega}\left(\frac{\theta_0}{\theta}-1\right)ds \leq \sum_{i=1}^{n-1}\int_\Omega r_i v_i dx. \qquad (2.7)$$

The bound for $\nabla\boldsymbol{v}$ comes from the positive definiteness of the *reduced* diffusion matrix $(M_{ij})_{i,j=1}^{n-1}$; see Lemma 6. Under suitable conditions on the heat conductivity and the reaction rates, this so-called entropy inequality provides gradient estimates for $\boldsymbol{v}$, $\log\theta$, $\theta$, and $\Pi\boldsymbol{q}$, but not for the full vector $\boldsymbol{q}$. Indeed, the relation $v_i = \log\rho_i - \log\rho_n$ can be inverted yielding

$$\rho_i = \frac{\rho^0\exp(v_i)}{\sum_{j=1}^n\exp(v_j)}, \quad i=1,\ldots,n-1, \quad \rho_n = \rho^0 - \sum_{j=1}^{n-1}\rho_j, \qquad (2.8)$$

which suggests to work with the *reduced* vector $\boldsymbol{\rho}' = (\rho_1,\ldots,\rho_{n-1})$. Moreover, this shows that $\rho_i$ stays bounded in some interval $(0,\rho^*)$ and, in view of the bound for $\nabla\boldsymbol{v}$, that $\nabla\boldsymbol{\rho}$ is bounded in $L^2(\Omega)$. Together with a bound for the (discrete) time derivative of $\rho_i$, we deduce the strong convergence of $\rho_i$ from the Aubin–Lions compactness lemma.

Still, there remains a difficulty. The estimate for $\kappa(\theta)^{1/2}\nabla\log\theta$ in $L^2(\Omega)$ from (2.7) is not sufficient to define $\kappa(\theta)\nabla\theta$ in the weak formulation. The idea is to derive better estimates for the temperature by using $\theta$ as a test function in the weak formulation of (1.2). If $\kappa(\theta) \geq c_\kappa\theta^2$ for some $c_\kappa > 0$ and $M_j/\theta$ is assumed to be bounded, then a formal computation, which is made precise in Lemma 7, gives

$$\frac{1}{2}\frac{d}{dt}\int_\Omega \rho^0\theta^2 dx + c_\kappa\int_\Omega\theta^2|\nabla\theta|^2 dx - \lambda\int_{\partial\Omega}(\theta_0-\theta)\theta ds \qquad (2.9)$$

$$= \sum_{j=1}^{n-1}\int_\Omega\frac{M_j}{\theta}\theta\nabla v_j\cdot\nabla\theta dx \leq \frac{c_\kappa}{2}\int_\Omega\theta^2|\nabla\theta|^2 dx + C\sum_{j=1}^{n-1}\int_\Omega|\nabla v_j|^2 dx.$$

Since $\nabla v_j$ is bounded in $L^2$, this yields uniform bounds for $\theta^2$ in $L^\infty(0,T;L^1(\Omega))$ and $L^2(0,T;H^1(\Omega))$. These estimates are sufficient to treat the term $\kappa(\theta)\nabla\theta$. The delicate point is to choose the approximate scheme in such a way that estimates (2.7) and (2.9) can be made rigorous; we refer to Section 2.4 for details.

## 2.3 Modeling

We consider an ideal fluid mixture consisting of $n$ components with the same molar masses in a fixed container $\Omega\subset\mathbb{R}^3$. The balance equations for the partial mass densities $\rho_i$ are given by

$$\partial_t\rho_i + \mathrm{div}(\rho_i v_i) = r_i, \quad i=1,\ldots,n,$$

where $v_i$ are the partial velocities and $r_i$ the reaction rates. Introducing the total mass density $\rho = \sum_{i=1}^n\rho_i$, the barycentric velocity $v = \rho^{-1}\sum_{i=1}^n\rho_i v_i$, and the diffusion fluxes

$J_i = \rho_i(v_i - v)$, we can reformulate the mass balances as

$$\partial_t \rho_i + \operatorname{div}(\rho_i v + J_i) = r_i, \quad i = 1, \ldots, n. \tag{2.10}$$

By definition, we have $\sum_{i=1}^n J_i = 0$, which means that the total mass density satisfies $\partial_t \rho + \operatorname{div}(\rho v) = 0$. We assume that the barycentric velocity vanishes, $v = 0$, i.e., the barycenter of the fluid is not moving. Consequently, the total mass density is constant in time.

The non-isothermal dynamics of the mixture is assumed to be given by the balance equations

$$\partial_t \rho_i + \operatorname{div} J_i = r_i, \quad \partial_t E + \operatorname{div} J_e = 0, \quad i = 1, \ldots, n,$$

where $J_e$ is the energy flux and $E$ the total energy. We suppose that the diffusion fluxes are proportional to the gradients of the thermo-chemical potentials $q_j$ and the temperature gradient (Soret effect) and that the energy flux is linear in the temperature gradient and the gradients of $q_j$ (Dufour effect):

$$J_i = -\sum_{j=1}^n M_{ij} \nabla q_j - M_i \nabla \frac{1}{\theta}, \quad i = 1, \ldots, n, \quad J_e = -\kappa(\theta) \nabla \theta - \sum_{j=1}^n M_j \nabla q_j.$$

The proportionality factor $\kappa(\theta)$ between the heat flux and the temperature gradient is the heat (or thermal) conductivity.

The thermo-chemical potentials and the total energy are determined in a thermodynamically consistent way from the free energy

$$\psi(\boldsymbol{\rho}, \theta) = \theta \sum_{i=1}^n \rho_i(\log \rho_i - 1) - \rho\theta(\log \theta - 1).$$

For simplicity, we have set the heat capacity equal to one. The physical entropy $s$, the chemical potentials $\mu_i$, and the total energy $E$ are defined by the free energy according to

$$s = -\frac{\partial \psi}{\partial \theta} = -\sum_{i=1}^n \rho_i(\log \rho_i - 1) + \rho \log \theta,$$

$$\mu_i = \frac{\partial \psi}{\partial \rho_i} = \theta(\log(\rho_i/\theta) + 1), \quad i = 1, \ldots, n,$$

$$E = \psi + \theta s = \rho\theta.$$

We introduce the mathematical entropy $h := -s$ and the thermo-chemical potentials $q_j = \mu_j/\theta = \log(\rho_j/\theta) + 1$ for $j = 1, \ldots, n$. These definitions lead to system (1.1)–(1.2). The Gibbs–Duhem relation yields the pressure $p = -\psi + \sum_{i=1}^n \rho_i \mu_i = \rho\theta$ of an ideal gas mixture. Note that we do not need a pressure blow-up at $\rho = 0$ to exclude vacuum or a superlinear growth in $\theta$ to control the temperature. Note also that, because of the nonvanishing pressure, one may criticize the choice of vanishing barycentric velocity. In the general case, the mass and energy balances need to be coupled with the momentum balance for $v$. Such systems, but only for isothermal or stationary systems, have been analyzed in, e.g., [BJPZ22, CJ15, DDGG20, Dru16]. The choice $v = 0$ is a mathematical simplification.

If the molar masses $m_i$ of the components are not the same, we need to modify the free energy according to [BJPZ22, Remark 1.2]

$$\psi = \theta \sum_{i=1}^{n} \frac{\rho_i}{m_i} \left( \log \frac{\rho_i}{m_i} - 1 \right) - c_W \rho \theta (\log \theta - 1),$$

where $c_W > 0$ is the heat capacity. For simplicity, we have set $m_i = 1$ and $c_W = 1$.

We show that the Maxwell–Stefan equations

$$\partial_t \rho_i + \operatorname{div} J_i = r_i, \quad d_i = -\sum_{j=1}^{n} b_{ij} \rho_i \rho_j \left( \frac{J_i}{\rho_i} - \frac{J_j}{\rho_j} \right), \quad i = 1, \dots, n, \tag{2.11}$$

with $b_{ij} = b_{ji} > 0$ can be formulated as (1.1) for a specific choice of $d_i$, $M_{ij}$, and $M_i$. The coefficients $b_{ij}$ may be interpreted as friction coefficients and can depend on $(\boldsymbol{\rho}, \theta)$; see [BD21, Section 4]. The equivalence between the Fick–Onsager and Maxwell-Stefan formulations was thoroughly investigated in [BD23], and we adapt their proof to our nonisothermal framework. For this, we introduce the matrix $B = (B_{ij})$ satisfying $B_{ii} = \sum_{j=1, j \neq i}^{n} b_{ij} \rho_j$ and $B_{ij} = -b_{ij} \rho_i$ for $j \neq i$. It is not invertible since $\boldsymbol{\rho} \in \ker(B)$, but its group inverse $B^{\#}$ exists uniquely, satisfying $BB^{\#} = B^{\#}B = I - (\boldsymbol{\rho}/\rho) \otimes \mathbf{1}$ and

$$\sum_{j=1}^{n} B_{ij}^{\#} \rho_j = 0, \quad \sum_{j=1}^{n} B_{ji}^{\#} = 0 \quad \text{for } i = 1, \dots, n. \tag{2.12}$$

Furthermore, we introduce the projection $P = (P_{ij}) = I - \mathbf{1} \otimes (\boldsymbol{\rho}/\rho)$ on $\operatorname{span}\{\boldsymbol{\rho}\}^{\perp}$.

**Proposition 5.** *Define the driving forces*

$$d_i = \rho_i \nabla \frac{\mu_i}{\theta} - \frac{\rho_i}{\rho \theta} \nabla (\rho \theta) - 2 \rho_i \theta \nabla \frac{1}{\theta} + q_i \rho_i \nabla \log \theta \quad \text{for } i = 1, \dots, n, \tag{2.13}$$

*where the numbers $q_i \in \mathbb{R}$ satisfy $\sum_{i=1}^{n} q_i \rho_i = 0$. Then (2.11) can be written as (1.1) with*

$$M_{ij} = \sum_{k=1}^{n} B_{ik}^{\#} \rho_k P_{kj}, \quad M_i = -\theta \sum_{k=1}^{n} B_{ik}^{\#} \rho_k q_k \quad \text{for } i, j = 1, \dots, n, \tag{2.14}$$

*where $(M_{ij})$ is symmetric and $M_{ij}$ and $M_i$ satisfy (1.5).*

The first three terms in the driving forces (2.13) are the same as [BD23, (4.18)] and [BD15, (2.11)], while the last term is motivated from [TA99, (A5)]. A computation shows that $\sum_{i=1}^{n} d_i = 0$ which is consistent with (2.11). It is argued in [BD23] that $M_{ij}$ is of the form $\rho_i(a_i(\boldsymbol{\rho}, \theta)\delta_{ij} + \rho_j S_{ij}(\boldsymbol{\rho}, \theta))$ for some functions $a_i$ and $S_{ij}$, and in the nondegenerate case, one may assume that $a_i(\boldsymbol{\rho}, \theta)$ stays positive when $\boldsymbol{\rho} \to \widetilde{\boldsymbol{\rho}}$ with $\widetilde{\rho}_i = 0$ [BD23, (6.6)]. This formulation motivates condition (2.6).

*Proof.* The proof is based on the equivalence between the Fick–Onsager and Maxwell–Stefan formulations elaborated in [BD23, Section 4] for the isothermal case. First, the driving forces can be formulated as

$$d_i = \rho_i \nabla \frac{\mu_i}{\theta} - \rho_i \nabla \log \frac{\rho}{\theta} + q_i \rho_i \nabla \log \theta,$$

which shows that

$$\sum_{j=1}^{n} \rho_j \nabla \frac{\mu_j}{\theta} = \sum_{j=1}^{n} \left( d_j + \rho_j \nabla \log \frac{\rho}{\theta} - q_j \rho_j \nabla \log \theta \right) = \rho \nabla \log \frac{\rho}{\theta}.$$

Consequently, another formulation is

$$d_i = \rho_i \nabla \frac{\mu_i}{\theta} - \frac{\rho_i}{\rho} \sum_{j=1}^{n} \rho_j \nabla \frac{\mu_j}{\theta} + q_i \rho_i \nabla \log \theta = \sum_{j=1}^{n} \rho_i P_{ij} \nabla \frac{\mu_j}{\theta} + q_i \rho_i \nabla \log \theta.$$

Setting $R = \operatorname{diag}(\rho_1, \ldots, \rho_n)$ and $\boldsymbol{q}^* = \operatorname{diag}(q_1 \rho_1, \ldots, q_n \rho_n)$, we obtain $\boldsymbol{d} = RP\nabla(\boldsymbol{\mu}/\theta) + \boldsymbol{q}^* \nabla \log \theta$. On the other hand, by (2.11),

$$d_i = -\left( \sum_{j=1, j\neq i}^{n} b_{ij} \rho_j \right) J_i + \sum_{j=1, j\neq i}^{n} b_{ij} \rho_i J_j = -\sum_{j=1}^{n} B_{ij} J_j.$$

This shows that $\boldsymbol{d} = -B\boldsymbol{J}$ and hence $\boldsymbol{J} = -B^{\#}\boldsymbol{d} = -B^{\#} RP\nabla(\boldsymbol{\mu}/\theta) - B^{\#}\boldsymbol{q}^* \nabla \log \theta$. Thus, defining $M_{ij}$ and $M_i$ as in (2.14), it follows that

$$J_i = -\sum_{j=1}^{n} M_{ij} \nabla \frac{\mu_j}{\theta} - M_i \nabla \frac{1}{\theta}.$$

The matrix $\tau = BR$ is symmetric and so does $\tau^{\#}$. Moreover, by [BD23, (4.26)], $B^{\#} = P^{\top} R \tau^{\#} P^{\top}$. Therefore, $M = B^{\#} RP = P^{\top} R \tau^{\#} RP$ is symmetric. We deduce from the properties (2.12) that

$$\sum_{j=1}^{n} M_{ij} = \sum_{j,k=1}^{n} B_{ik}^{\#} \rho_k \left( \delta_{kj} - \frac{\rho_j}{\rho} \right) = 0, \quad \sum_{i=1}^{n} M_i = -\theta \sum_{j=1}^{n} \left( \sum_{i=1}^{n} B_{ij}^{\#} \right) \rho_j q_j = 0.$$

This finishes the proof. $\qquad \square$

## 2.4 Proof of Existence – Nondegenerate Case (Theorem 3)

The idea of the proof is to reformulate equations (1.1)–(1.2) in terms of the relative potentials $v_i$, to approximate the resulting equations by an implicit Euler scheme, and to add some higher-order regularizations in space for the variables $v_i$ and $w = \log \theta$. The de-regularization limit is based on the compactness coming from the entropy estimates and an estimate for the temperature.

Set $w_0 = \log \theta_0$, $\varepsilon > 0$, $N \in \mathbb{N}$, and $\tau = T/N > 0$. To simplify the notation, we set $\boldsymbol{v} = (\boldsymbol{v}', 0) = (v_1, \ldots, v_{n-1}, 0)$ and $\bar{\boldsymbol{v}} = (\bar{v}_1, \ldots, \bar{v}_{n-1}, 0)$. Let $(\bar{\boldsymbol{v}}, \bar{w}) \in L^{\infty}(\Omega; \mathbb{R}^{n+1})$ be given, and set $\rho_i(\boldsymbol{v}) = \rho^0 e^{v_i} / \sum_{j=1}^{n} e^{v_j}$ for $i = 1, \ldots, n-1$, $\rho_n = \rho^0 - \sum_{i=1}^{n-1} \rho_i$, and $q_i = \log \rho_i - w$ for $i = 1, \ldots, n$. We define the approximate scheme

$$0 = \frac{1}{\tau} \int_{\Omega} (\rho_i(\boldsymbol{v}) - \bar{\rho}_i(\bar{\boldsymbol{v}})) \phi_i dx + \int_{\Omega} \left( \sum_{j=1}^{n-1} M_{ij}(\boldsymbol{\rho}, e^w) \nabla v_j - M_i(\boldsymbol{\rho}, e^w) e^{-w} \nabla w \right) \cdot \nabla \phi_i dx$$

(2.15)

$$+ \varepsilon \int_\Omega \big(D^2 v_i : D^2 \phi_i + v_i \phi_i\big)dx - \int_\Omega r_i(\Pi \boldsymbol{q}, e^w)\phi_i dx,$$

$$0 = \frac{1}{\tau} \int_\Omega (E - \bar{E})\phi_0 dx + \int_\Omega \Big(\kappa(\theta)\nabla\theta + \sum_{j=1}^{n-1} M_j(\boldsymbol{\rho}, e^w)\nabla v_j\Big) \cdot \nabla\phi_0 dx \qquad (2.16)$$

$$- \lambda \int_{\partial\Omega} (\theta_0 - \theta)\phi_0 ds + \varepsilon \int_\Omega e^w \big(D^2 w : D^2 \phi_0 + |\nabla w|^2 \nabla w \cdot \nabla\phi_0\big)dx$$

$$+ \varepsilon \int_\Omega (e^{w_0} + e^w)(w - w_0)\phi_0 dx$$

for test functions $\phi_i \in H^2(\Omega)$, $i = 0, \ldots, n-1$. Here, $D^2 u$ is the Hessian matrix of the function $u$, ":" denotes the Frobenius matrix product, and $E = \rho^0 \theta$, $\bar{E} = \rho^0 \bar{\theta}$. The lower-order regularization $\varepsilon(e^{w_0} + e^w)(w - w_0)$ yields an $L^2(\Omega)$ estimate for $w$. Furthermore, the higher-order regularization guarantees that $v_i, w \in H^2(\Omega) \hookrightarrow L^\infty(\Omega)$, while the $W^{1,4}(\Omega)$ regularization term for $w$ allows us to estimate the higher-order terms when using the test function $e^{-w_0} - e^{-w}$.

*Step 1: solution of the linearized approximate problem.* In order to define the fixed-point operator, we need to solve a linearized problem. To this end, let $y^* = (\boldsymbol{v}^*, w^*) \in W^{1,4}(\Omega; \mathbb{R}^n)$ and $\sigma \in [0, 1]$ be given. We want to find the unique solution $y = (\boldsymbol{v}', w) \in H^2(\Omega; \mathbb{R}^n)$ to the linear problem

$$a(y, \phi) = \sigma F(\phi) \quad \text{for all } \phi = (\phi_0, \ldots, \phi_{n-1}) \in H^2(\Omega; \mathbb{R}^n), \qquad (2.17)$$

where

$$a(y, \phi) = \int_\Omega \sum_{i,j=1}^{n-1} M_{ij}(\boldsymbol{\rho}^*, e^{w^*})\nabla v_j \cdot \nabla\phi_i dx + \int_\Omega \kappa(e^{w^*})e^{w^*}\nabla w \cdot \nabla\phi_0 dx$$

$$+ \varepsilon \int_\Omega \sum_{i=1}^{n-1} \big(D^2 v_i : D^2 \phi_i + v_i \phi_i\big)$$

$$+ \varepsilon \int_\Omega e^{w^*}\big(D^2 w : D^2 \phi_0 + |\nabla w^*|^2 \nabla w \cdot \nabla\phi_0\big) + \varepsilon \int_\Omega (e^{w_0} + e^{w^*})w\phi_0 dx,$$

$$F(\phi) = -\frac{1}{\tau} \int_\Omega \sum_{i=1}^{n-1} (\rho_i^* - \bar{\rho}_i)\phi_i dx - \frac{1}{\tau} \int_\Omega (E^* - \bar{E})\phi_0 dx + \lambda \int_{\partial\Omega} (e^{w_0} - e^{w^*})\phi_0 dx$$

$$+ \int_\Omega \sum_{i=1}^{n-1} M_i(\boldsymbol{\rho}^*, e^{w^*})e^{-w^*}\nabla w^* \cdot \nabla\phi_i dx - \int_\Omega \sum_{j=1}^{n-1} M_j(\boldsymbol{\rho}^*, e^{w^*})\nabla v_j^* \cdot \nabla\phi_0 dx$$

$$+ \int_\Omega \sum_{i=1}^{n} r_i(\Pi \boldsymbol{q}^*, e^{w^*})\phi_i dx + \varepsilon \int_\Omega (e^{w_0} + e^{w^*})w_0\phi_0 dx$$

and $\rho_i^* = \rho_i(\boldsymbol{v}^*)$, $\rho^* = \sum_{i=1}^{n} \rho_i^*$, $E^* = \rho^0 e^{w^*}$. By Hypothesis (H3) and the generalized Poincaré inequality [Tem97, Chap. 2, Sec. 1.4], we have

$$a(y, y) \geq \varepsilon \int_\Omega \big(|D^2 \boldsymbol{v}|^2 + |\boldsymbol{v}|^2\big)dx + \varepsilon \int_\Omega e^{w^*}(|D^2 w|^2 + w^2)dx \geq \varepsilon C\big(\|\boldsymbol{v}\|_{H^2(\Omega)}^2 + \|w\|_{H^2(\Omega)}^2\big).$$

Thus, $a$ is coercive. Moreover, $a$ and $F$ are continuous on $H^2(\Omega; \mathbb{R}^n)$. The Lax–Milgram lemma shows that (2.17) possesses a unique solution $(\boldsymbol{v}', w) \in H^2(\Omega; \mathbb{R}^n)$.

*Step 2: solution of the approximate problem.* The previous step shows that the fixed-point operator $S : W^{1,4}(\Omega; \mathbb{R}^n) \times [0,1] \to W^{1,4}(\Omega; \mathbb{R}^n)$, $S(y^*, \sigma) = y$, where $y = (\boldsymbol{v}', w)$ solves (2.17), is well defined. It holds that $S(y, 0) = 0$, $S$ is continuous, and since $S$ maps to $H^2(\Omega; \mathbb{R}^n)$, which is compactly embedded into $W^{1,4}(\Omega; \mathbb{R}^n)$, it is also compact. It remains to determine a uniform bound for all fixed points $y$ of $S(\cdot, \sigma)$, where $\sigma \in [0,1]$. Let $y$ be such a fixed point. Then $y \in H^2(\Omega; \mathbb{R}^n)$ solves (2.17) with $(\boldsymbol{v}^*, w^*)$ replaced by $y = (\boldsymbol{v}', w)$. With the test functions $\phi_i = v_i$ for $i = 1, \ldots, n-1$ and $\phi_0 = e^{-w_0} - e^{-w}$ (we need this test function since $\phi_0 = -e^{-w}$ does not allow us to control the lower-order term), we obtain

$$
0 = \frac{\sigma}{\tau} \int_\Omega \sum_{i=1}^{n-1} (\rho_i(\boldsymbol{v}) - \rho_i(\bar{\boldsymbol{v}})) v_i dx + \frac{\sigma}{\tau} \int_\Omega (E - \bar{E})(-e^{-w}) dx + \frac{\sigma}{\tau} \int_\Omega (E - \bar{E}) e^{-w_0} dx
$$
$$
+ \int_\Omega \sum_{i,j=1}^{n-1} M_{ij} \nabla v_i \cdot \nabla v_j dx + \int_\Omega \kappa(e^w) e^w \nabla w \cdot \nabla(-e^{-w}) dx - \sigma \int_\Omega \sum_{i=1}^{n-1} r_i v_i dx
$$
$$
- \sigma \int_\Omega \sum_{j=1}^{n-1} M_j e^{-w} \nabla w \cdot \nabla v_j dx + \sigma \int_\Omega \sum_{j=1}^{n-1} M_j \nabla v_j \cdot \nabla(-e^{-w}) dx
$$
$$
- \sigma \lambda \int_{\partial\Omega} (e^{w_0} - e^w)(e^{-w_0} - e^{-w}) dx + \varepsilon \int_\Omega \sum_{i=1}^{n-1} (|D^2 v_i|^2 + v_i^2) dx
$$
$$
+ \varepsilon \int_\Omega (e^{w_0} + e^w)(w - \sigma w_0)(e^{-w_0} - e^{-w}) dx
$$
$$
+ \varepsilon \int_\Omega (|D^2 w|^2 - D^2 w : \nabla w \otimes \nabla w + |\nabla w|^4) dx
$$
$$
=: I_1 + \cdots + I_{12}.
$$

We see immediately that $I_7 + I_8 = 0$. Furthermore,

$$
I_1 + I_2 = \frac{\sigma}{\tau} \int_\Omega \left( \sum_{i=1}^{n-1} (\rho_i - \bar{\rho}_i) \frac{\partial h}{\partial \rho_i} + (\theta - \bar{\theta}) \frac{\partial h}{\partial \theta} \right) dx.
$$

The function $(\boldsymbol{\rho}', \theta) \mapsto h(\boldsymbol{\rho}', \theta) = \sum_{i=1}^n \rho_i(\log \rho_i - 1) - \rho^0 \log \theta$ with $\rho_n = \rho^0 - \sum_{i=1}^{n-1} \rho_i$ is convex, since the second derivatives are given by

$$
\frac{\partial^2 h}{\partial \rho_i^2} = \frac{1}{\rho_i} + \frac{1}{\rho_n}, \quad \frac{\partial^2 h}{\partial \theta^2} = \frac{\rho^0}{\theta^2}, \quad \frac{\partial^2 h}{\partial \rho_i \partial \theta} = 0, \quad \frac{\partial^2 h}{\partial \rho_i \partial \rho_j} = \frac{1}{\rho_n},
$$

hence we can conclude in the same way as in [JS13] that the Hessian is positive definite by Sylvester's criterion. This shows that

$$
h(\boldsymbol{\rho}', \theta) - h(\bar{\boldsymbol{\rho}}', \bar{\theta}) \leq \sum_{i=1}^{n-1} \frac{\partial h}{\partial \rho_i}(\boldsymbol{\rho}', \theta)(\rho_i - \bar{\rho}_i) + \frac{\partial h}{\partial \theta}(\boldsymbol{\rho}', \theta)(\theta - \bar{\theta})
$$

and consequently,

$$I_1 + I_2 \geq \frac{\sigma}{\tau} \int_\Omega \left( h(\boldsymbol{\rho}', \theta) - h(\bar{\boldsymbol{\rho}}', \bar{\theta}) \right) dx.$$

For the estimate of $I_4$, we need the following lemma.

**Lemma 6.** *Let the matrix $(M_{ij}) \in \mathbb{R}^{n \times n}$ satisfy* (1.5) *and* (1.6). *Then*

$$\sum_{i,j=1}^{n-1} M_{ij}(z_i - z_n)(z_j - z_n) \geq \frac{c_M}{n} \sum_{i=1}^{n-1} |z_i - z_n|^2.$$

*Proof.* We use (1.5) and then (1.6) to find for any $\boldsymbol{z} \in \mathbb{R}^n$ that

$$\sum_{i,j=1}^{n-1} M_{ij}(z_i - z_n)(z_j - z_n) = \sum_{i,j=1}^{n} M_{ij} z_i z_j \geq c_M |\Pi \boldsymbol{z}|^2. \tag{2.18}$$

Therefore, we want to prove

$$c_M |\Pi z|^2 \geq \frac{c_M}{n} \sum_{i=1}^{n} (z_i - z_n)^2. \tag{2.19}$$

To prove this, we first recall

$$|\Pi z|^2 = \sum_{i=1}^{n} \left( z_i - \frac{1}{n} \sum_{j=1}^{n} z_j \right)^2.$$

We observe with $\bar{z} = \frac{1}{n} \sum_{j=1}^{n} z_j$ that the inequality we want to prove is equivalent to

$$\sum_{i=1}^{n} \left( z_i - \frac{1}{n} \sum_{j=1}^{n} z_j \right)^2 - \frac{1}{n} \sum_{i=1}^{n-1} (z_i - z_n)^2 \geq 0$$

$$\Leftrightarrow \sum_{i=1}^{n} \left( z_i^2 - 2 z_i \bar{z} + \bar{z}^2 \right) - \frac{1}{n} \sum_{i=1}^{n} \left( z_i^2 - 2 z_n z_i + z_n^2 \right) \geq 0$$

$$\Leftrightarrow \sum_{i=1}^{n} z_i^2 - n \bar{z}^2 - \frac{1}{n} \sum_{i=1}^{n} z_i^2 + 2 \frac{1}{n} \sum_{i=1}^{n} z_n z_i - z_n^2 \geq 0$$

$$\Leftrightarrow \frac{n-1}{n} \sum_{i=1}^{n} z_i^2 - \frac{1}{n} \left( \sum_{i=1}^{n} z_i \right)^2 + 2 \frac{1}{n} \sum_{i=1}^{n} z_n z_i - z_n^2 \geq 0$$

$$\Leftrightarrow \frac{n-1}{n} \sum_{i=1}^{n-1} z_i^2 - \frac{1}{n} \left( \sum_{i=1}^{n} z_i - z_n \right)^2 \geq 0$$

$$\Leftrightarrow \sum_{i=1}^{n-1} z_i^2 \geq \frac{1}{n-1} \left( \sum_{i=1}^{n-1} z_i \right)^2.$$

Since $x \mapsto x^2$ is a convex function, Jensen inequality implies

$$(n-1)\left(\frac{\sum_{i=1}^{n-1} z_i}{n-1}\right)^2 \leq (n-1)\frac{\sum_{i=1}^{n-1} z_i^2}{n-1} = \sum_{i=1}^{n-1} z_i^2,$$

we conclude that (2.19) holds.

$\square$

By Lemma 6 and Hypothesis (H5),

$$I_4 = \frac{1}{2} \int_\Omega \left( \sum_{i,j=1}^{n} M_{ij} \nabla q_i \cdot \nabla q_j + \sum_{i,j=1}^{n-1} M_{ij} \nabla v_i \cdot \nabla v_j \right) dx$$

$$\geq \frac{c_M}{2} \int_\Omega |\nabla \Pi \boldsymbol{q}|^2 dx + \frac{c_M}{2n} \int_\Omega |\nabla \boldsymbol{v}|^2 dx,$$

$$I_6 = \sigma \int_\Omega \sum_{i=1}^{n} r_i q_i dx \geq \sigma c_r \int_\Omega |\Pi \boldsymbol{q}|^2 dx.$$

Furthermore, we observe

$$I_{11} = \varepsilon \int_\Omega (e^{w_0} + e^w)(w - \sigma w_0)(e^{-w_0} - e^{-w}) dx$$

$$= 2\varepsilon \int_\Omega (w - \sigma w_0) \sinh(w - w_0) dx.$$

Rewriting the integrand and using $\cosh(x) \geq x$ for all $x \in \mathbb{R}$, we find

$$(w - \sigma w_0)\sinh(w - w_0) = (w - w_0)(w - \sigma w_0)\int_0^1 \cosh(s(w - w_0)) ds$$

$$= (w^2 - \sigma w_0 w - w_0 w + \sigma w_0^2)\int_0^1 \cosh(s(w - w_0)) ds$$

$$\geq \frac{1}{2}w^2 + \left(\frac{1}{2}w^2 - (\sigma + 1)w_0 w\right)\frac{\sinh(w - w_0)}{w - w_0}.$$

Since $\sinh(w - w_0)/(w - w_0) \geq 1$, the lower bound of the right hand side depends on the sign of $1/2w^2 - (\sigma + 1)w_0 w$. Since $\sigma$ and $w_0$ are nonnegative, we find

$$\frac{1}{2}w^2 - (\sigma + 1)w_0 w < 0 \Rightarrow 0 < w \leq 2(\sigma + 1)w_0.$$

Therefore, it is sufficient to prove that the function

$$g \colon [0, 4w_0] \times [0, 1] \to \mathbb{R}$$

$$(w, \sigma) \mapsto \left(\frac{1}{2}w^2 - (\sigma + 1)w_0 w\right)\frac{\sinh(w - w_0)}{w - w_0}$$

attains a minimum. However, $g$ is a continuous function on a compact set and therefore has a minimum, i.e. it exists $(\hat{w}, \hat{\sigma})$ such that $g(\hat{w}, \hat{\sigma}) \leq g(w, \sigma)$ for all $(w, \sigma) \in [0, 4w_0] \times [0, 1]$. Hence, we conclude

$$I_{11} \geq \frac{1}{2} w^2 + g(\hat{w}, \hat{\sigma}).$$

Next, we have

$$I_5 = \int_\Omega \kappa(e^w) |\nabla w|^2 dx,$$

$$I_9 = 2\sigma\lambda \int_{\partial\Omega} (\cosh(w_0 - w) - 1) ds \geq 0,$$

$$I_{12} = \frac{\varepsilon}{2} \int_\Omega \left( |D^2 w|^2 + |D^2 w - \nabla w \otimes \nabla w|^2 + |\nabla w|^4 \right) dx.$$

Summarizing these estimates and applying the generalized Poincaré inequality, we arrive at the *discrete entropy inequality*

$$
\frac{\sigma}{\tau} \int_\Omega \left( h(\boldsymbol{\rho}', \theta) + e^{-w_0} E \right) dx + \frac{c_M}{2} \int_\Omega \left( \frac{1}{n} |\nabla \boldsymbol{v}|^2 + |\nabla \Pi \boldsymbol{q}|^2 + \sigma c_r |\Pi \boldsymbol{q}|^2 \right) dx
$$

$$
+ \varepsilon C \left( \|\boldsymbol{v}\|^2_{H^2(\Omega)} + \|w\|^2_{H^2(\Omega)} + \|w\|^4_{W^{1,4}(\Omega)} \right) + \int_\Omega \kappa(e^w) |\nabla w|^2 dx
$$

$$
\leq \frac{\sigma}{\tau} \int_\Omega \left( h(\bar{\boldsymbol{\rho}}', \bar{\theta}) + e^{-w_0} \bar{E} \right) dx + 2\varepsilon \|w_0\|^2_{L^2(\Omega)} + C(w_0), \tag{2.20}
$$

where $C(w_0) > 0$ is the constant, coming from the estimate of $I_{11}$. We observe that the left-hand side is bounded from below since $-\rho^0 \log\theta + e^{-w_0} E = \rho^0(-\log\theta + \theta/\theta_0)$ is bounded from below. The bound for $\Pi \boldsymbol{q}$ implies an $L^2(\Omega)$ bound for $\boldsymbol{v}$ since $|\boldsymbol{v}|^2 \leq n|\Pi \boldsymbol{q}|^2$; see the proof of Lemma 6.

Estimate (2.20) gives a uniform bound for $(\boldsymbol{v}', w)$ in $H^2(\Omega; \mathbb{R}^n)$ and consequently also in $W^{1,4}(\Omega; \mathbb{R}^n)$, which proves the claim. We infer from the Leray–Schauder fixed-point theorem that there exists a solution $(\boldsymbol{v}', w)$ to (2.15)–(2.16).

*Step 3: temperature estimate.* We need a better estimate for the temperature.

**Lemma 7.** *Let $(\boldsymbol{\rho}, w)$ be a solution to (2.15)–(2.16) and set $\theta = e^w$. Then there exists a constant $C > 0$ independent of $\varepsilon$ and $\tau$ such that*

$$\frac{1}{\tau} \int_\Omega \rho^0 \theta^2 dx + \frac{1}{2} \int_\Omega \kappa(\theta) |\nabla\theta|^2 dx \leq C + \frac{1}{\tau} \int_\Omega \rho^0 \bar{\theta}^2 dx + C \int_\Omega |\nabla\boldsymbol{v}|^2 dx.$$

*Proof.* We use the test function $\theta$ in (2.16). Observing that $(E - \bar{E})\theta = \rho^0(\theta - \bar{\theta})\theta \geq (\rho^0/2)(\theta^2 - \bar{\theta}^2)$ and that $\kappa(\theta) \geq c_\kappa(1 + \theta^2)$ by Hypothesis (H4), we find that

$$
\frac{1}{2\tau} \int_\Omega \rho^0(\theta^2 - \bar{\theta}^2) dx + \frac{1}{2} \int_\Omega \kappa(\theta)|\nabla\theta|^2 dx + \frac{c_\kappa}{2} \int_\Omega \theta^2 |\nabla\theta|^2 dx - \lambda \int_{\partial\Omega} (\theta_0 - \theta)\theta dx
$$

$$
\leq -\sum_{j=1}^{n-1} \int_\Omega M_j \nabla v_j \cdot \nabla\theta dx - \varepsilon \int_\Omega \theta(D^2 \log\theta : D^2\theta + |\nabla \log\theta|^2 \nabla \log\theta \cdot \nabla\theta) dx
$$

$$- \varepsilon \int_\Omega (\theta_0 + \theta)(\log \theta - \log \theta_0)\theta dx$$
$$=: J_1 + J_2 + J_3. \tag{2.21}$$

Since $M_j/\theta$ is assumed to be bounded,

$$J_1 \leq \frac{c_\kappa}{2} \int_\Omega \theta^2 |\nabla \theta|^2 dx + C \sum_{j=1}^{n-1} \int_\Omega |\nabla v_j|^2 dx.$$

Furthermore,

$$J_2 = -\varepsilon \int_\Omega \left( -\frac{1}{\theta} \nabla \theta \cdot D^2 \theta \nabla \theta + |D^2 \theta|^2 + \frac{1}{\theta^2} |\nabla \theta|^4 \right) dx$$
$$= -\frac{\varepsilon}{2} \int_\Omega \left( |D^2 \theta|^2 + \frac{1}{\theta^2} |\nabla \theta|^4 + \left| D^2 \theta - \frac{1}{\theta} \nabla \theta \otimes \nabla \theta \right|^2 \right) dx \leq 0.$$

The last integral $J_3$ is bounded since $-\theta^2 \log \theta$ is the dominant term. The last term on the left-hand side of (2.21) is bounded from below by $-(\lambda/2)\int_{\partial\Omega} \theta_0^2 dx$, which finishes the proof. □

**Remark 8.** *Better estimates can be derived if we assume that $\kappa(\theta) \geq c_\kappa(1 + \theta^{\alpha+1})$ for $\alpha \in (1, 2)$. Indeed, using $\theta^\alpha$ as a test function in (2.16), we find that*

$$\frac{1}{\tau} \int_\Omega \rho^0 (\theta - \bar{\theta})\theta^\alpha dx + \alpha c_\kappa \int_\Omega \theta^{2\alpha} |\nabla \theta|^2 dx - \lambda \int_{\partial\Omega} (\theta_0 - \theta)\theta^\alpha dx$$
$$\leq -\alpha \sum_{j=1}^{n-1} \int_\Omega M_j \theta^{\alpha-1} \nabla v_j \cdot \nabla \theta dx - \varepsilon \int_\Omega (\theta_0 + \theta)(\log \theta - \log \theta_0)\theta^\alpha dx$$
$$- \varepsilon \int_\Omega \theta \left( D^2 \log \theta : D^2 \theta^\alpha + |\nabla \log \theta|^2 \nabla \log \theta \cdot \nabla \theta^\alpha \right) dx$$
$$=: J_4 + J_5 + J_6. \tag{2.22}$$

*A tedious but straightforward computation shows that $J_6 \geq 0$ if $\alpha \in (1, 2)$. Furthermore, since $M_j/\theta$ is bounded,*

$$J_4 \leq \frac{\alpha c_\kappa}{2} \int_\Omega \theta^{2\alpha} |\nabla \theta|^2 dx + C \sum_{j=1}^{n-1} \int_\Omega |\nabla v_j|^2 dx.$$

*The first integral on the right-hand side is controlled by the left-hand side of (2.22). This yields a bound for $\theta^{\alpha+1} \in L^\infty(0, T; L^1(\Omega)) \cap L^2(0, T; H^1(\Omega)) \subset L^{8/3}(\Omega_T)$ (see Lemma 10) and consequently $\theta \in L^{8(\alpha+1)/3}(\Omega_T)$, which is better than the result in Lemma 10.* □

*Step 4: uniform estimates.* Let $((\boldsymbol{v}')^k, w^k)$ be a solution to (2.15)–(2.16) for given $(\boldsymbol{v}')^{k-1} = \bar{\boldsymbol{v}}'$ and $w^{k-1} = \bar{w}$, where $k \in \mathbb{N}$. We set

$$\theta^k = \exp(w^k), \quad \rho_i^k = \exp(w^k + q_i^k) = \frac{\rho^0 e^{v_i^k}}{\sum_{j=1}^n e^{v_j^k}}$$

for $i = 1, \ldots, n-1$, and $E^k = \rho^0 \theta^k$. We introduce piecewise constant functions in time. For this, let $\rho_i^{(\tau)}(x,t) = \rho_i^k(x)$, $\theta^{(\tau)}(x,t) = \theta^k(x)$, $v_i^{(\tau)}(x,t) = v_i^k(x)$, $q_i^{(\tau)} = \log(\rho_i^{(\tau)}/\theta^{(\tau)})$, and $E^{(\tau)}(x,t) = E^k(x)$ for $x \in \Omega$, $t \in ((k-1)\tau, k\tau]$, $k = 1, \ldots, N$. At time $t = 0$, we set $\rho_i^{(\tau)}(x,0) = \rho_i^0(x)$ and $\theta^{(\tau)}(x,0) = \theta^0(x)$ for $x \in \Omega$. Furthermore, we introduce the shift operator $(\sigma_\tau \rho_i^{(\tau)})(x,t) = \rho_i^{k-1}(x)$ for $x \in \Omega$, $t \in ((k-1)\tau, k\tau]$. Let $(\boldsymbol{\rho}')^{(\tau)} = (\rho_1^{(\tau)}, \ldots, \rho_{n-1}^{(\tau)})$. Then $((\boldsymbol{\rho}')^{(\tau)}, \theta^{(\tau)})$ solves (see (2.15)–(2.16))

$$0 = \frac{1}{\tau} \int_0^T \int_\Omega (\rho_i^{(\tau)} - \sigma_\tau \rho_i^{(\tau)}) \phi_i \, dx \, dt \tag{2.23}$$

$$+ \int_0^T \int_\Omega \left( \sum_{j=1}^{n-1} M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)} + M_i(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla \frac{1}{\theta^{(\tau)}} \right) \cdot \nabla \phi_i \, dx \, dt$$

$$+ \varepsilon \int_0^T \int_\Omega \left( D^2 v_i^{(\tau)} : D^2 \phi_i + v_i^{(\tau)} \phi_i \right) dx \, dt - \int_0^T \int_\Omega r_i(\Pi \boldsymbol{q}^{(\tau)}, \theta^{(\tau)}) \phi_i \, dx \, dt,$$

$$0 = \frac{1}{\tau} \int_0^T \int_\Omega (E^{(\tau)} - \sigma_\tau E^{(\tau)}) \phi_0 \, dx \, dt - \lambda \int_0^T \int_{\partial \Omega} (\theta_0 - \theta^{(\tau)}) \phi_0 \, ds \, dt \tag{2.24}$$

$$+ \int_0^T \int_\Omega \left( \kappa(\theta^{(\tau)}) \nabla \theta^{(\tau)} + \sum_{j=1}^{n-1} M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)} \right) \cdot \nabla \phi_0 \, dx \, dt$$

$$+ \varepsilon \int_0^T \int_\Omega \theta^{(\tau)} \left( D^2 \log \theta^{(\tau)} : D^2 \phi_0 + |\nabla \log \theta^{(\tau)}|^2 \nabla \log \theta^{(\tau)} \cdot \nabla \phi_0 \right) dx \, dt$$

$$+ \varepsilon \int_0^T \int_\Omega (\theta_0 + \theta^{(\tau)})(\log \theta^{(\tau)} - \log \theta_0) \phi_0 \, dx \, dt.$$

The discrete entropy inequality (2.20) and the $L^\infty$ bound for $\rho_i^{(\tau)}$ imply the following uniform bounds:

$$\|\rho_i^{(\tau)}\|_{L^\infty(0,T;L^\infty(\Omega))} + \|\theta^{(\tau)}\|_{L^\infty(0,T;L^1(\Omega))} \leq C,$$

$$\|v_i^{(\tau)}\|_{L^2(0,T;H^1(\Omega))} + \|\kappa(\theta^{(\tau)})^{1/2} \nabla \log \theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C,$$

$$\varepsilon^{1/2} \|v_i^{(\tau)}\|_{L^2(0,T;H^2(\Omega))} + \varepsilon^{1/2} \|\log \theta^{(\tau)}\|_{L^2(0,T;H^2(\Omega))} \leq C,$$

$$\varepsilon^{1/4} \|\log \theta^{(\tau)}\|_{L^4(0,T;W^{1,4}(\Omega))} \leq C,$$

for all $i = 1, \ldots, n-1$, where $C > 0$ is independent of $\varepsilon$ and $\tau$. Hypothesis (H4) yields

$$\|\nabla \theta^{(\tau)}\|_{L^2(\Omega_T)} + \|\nabla \log \theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C. \tag{2.25}$$

**Lemma 9** (Estimates for the temperature). *There exists a constant $C > 0$ which does not depend on $\varepsilon$ or $\tau$ such that*

$$\|\theta^{(\tau)}\|_{L^2(0,T;H^1(\Omega))} + \|\log \theta^{(\tau)}\|_{L^2(0,T;H^1(\Omega))} \leq C. \tag{2.26}$$

*Proof.* The entropy inequality shows that $-\log \theta^{(\tau)} + \theta^{(\tau)}$ is uniformly bounded from above, which shows that $|\log \theta^{(\tau)}|$ is uniformly bounded too and hence, $\log \theta^{(\tau)}$ is bounded in

$L^\infty(0,T;L^1(\Omega))$. Together with the $L^\infty(0,T;L^1(\Omega))$ bound for $\theta^{(\tau)}$, estimate (2.25), and the Poincaré–Wirtinger inequality, we find that

$$\|\theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C\|\theta^{(\tau)}\|_{L^2(0,T;L^1(\Omega))} + \|\nabla\theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C,$$
$$\|\log\theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C\|\log\theta^{(\tau)}\|_{L^2(0,T;L^1(\Omega))} + \|\nabla\log\theta^{(\tau)}\|_{L^2(\Omega_T)} \leq C,$$

from which we conclude the proof. □

We proceed by proving more uniform estimates. Because of the $L^2(\Omega_T)$ bound of $\nabla v_i^{(\tau)}$ and

$$\int_0^T\int_\Omega |\nabla\rho_i^{(\tau)}|^2 dxdt = \int_0^T\int_\Omega |\nabla\rho^0|^2 \left|\frac{\exp(v_i^{(\tau)})}{\sum_{j=1}^n \exp(v_j^{(\tau)})}\right|^2 dxdt$$
$$+ \int_0^T\int_\Omega \left|\frac{\exp(v_i^{(\tau)})\nabla v_i^{(\tau)}}{\sum_{j=1}^n \exp(v_j^{(\tau)})} - \frac{\exp(v_i^{(\tau)})\sum_{j=1}^n \exp(v_j^{(\tau)})\nabla v_j^{(\tau)}}{(\sum_{j=1}^n \exp(v_j^{(\tau)}))^2}\right|^2 dxdt$$
$$\leq \int_0^T\int_\Omega |\nabla\rho^0|^2 dxdt + 2\int_0^T\int_\Omega |\nabla\boldsymbol{v}|^2 dxdt \leq C, \tag{2.27}$$

$(\nabla\rho_i^{(\tau)})$ is bounded in $L^2(\Omega_T)$ and, taking into account the $L^\infty$ bound for $\rho_i^{(\tau)}$, the family $(\rho_i^{(\tau)})$ is bounded in $L^2(0,T;H^1(\Omega))$. By Lemma 7 and Hypothesis (H4), $(\nabla(\theta^{(\tau)})^2)$ is bounded in $L^2(\Omega_T)$. Therefore, since $((\theta^{(\tau)})^2)$ is bounded in $L^1(\Omega_T)$, the Poincaré–Wirtinger inequality gives a uniform bound for $(\theta^{(\tau)})^2$ in $L^2(0,T;H^1(\Omega))$. These bounds yields higher integrability of $\theta^{(\tau)}$, as shown in the following lemma.

**Lemma 10.** *There exists $C > 0$ independent of $\varepsilon$ and $\tau$ such that $(\theta^{(\tau)})$ is bounded in $L^{16/3}(\Omega_T)$.*

*Proof.* We deduce from the bound for $(\theta^{(\tau)})^2$ in $L^2(0,T;H^1(\Omega)) \subset L^2(0,T;L^6(\Omega))$ that $(\theta^{(\tau)})$ is bounded in $L^4(0,T;L^{12}(\Omega))$. By interpolation with $1/r = \alpha/12 + (1-\alpha)/2$ and $r\alpha = 4$,

$$\|\theta^{(\tau)}\|_{L^r(\Omega_T)}^r = \int_0^T \|\theta^{(\tau)}\|_{L^r(\Omega)}^r dt \leq \int_0^T \|\theta^{(\tau)}\|_{L^{12}(\Omega)}^{r\alpha}\|\theta^{(\tau)}\|_{L^2(\Omega)}^{r(1-\alpha)} dt$$
$$\leq \|\theta^{(\tau)}\|_{L^\infty(0,T;L^2(\Omega))}^{r(1-\alpha)}\int_0^T \|\theta^{(\tau)}\|_{L^{12}(\Omega)}^4 dt \leq C.$$

The solution of $1/r = \alpha/12 + (1-\alpha)/2$ and $r\alpha = 4$ is $\alpha = 3/4$ and $r = 16/3$. □

**Lemma 11.** *There exists $C > 0$ independent of $\varepsilon$ and $\tau$ such that*

$$\tau^{-1}\|\rho_i^{(\tau)} - \sigma_\tau\rho_i^{(\tau)}\|_{L^2(0,T;H^2(\Omega)')} + \tau^{-1}\|\theta^{(\tau)} - \sigma_\tau\theta^{(\tau)}\|_{L^{16/15}(0,T;W^{2,16}(\Omega)')} \leq C. \tag{2.28}$$

*Proof.* Let $\phi_0 \in L^{16}(0,T;W^{2,16}(\Omega))$, $\phi_1,\dots,\phi_{n-1} \in L^2(0,T;H^2(\Omega))$ and set $M_i^{(\tau)} = M_i(\boldsymbol{\rho}^{(\tau)},\theta^{(\tau)})$, $r_i^{(\tau)} = r_i(\boldsymbol{\rho}^{(\tau)},\theta^{(\tau)})$ for $i = 1,\dots,n-1$. It follows from (2.23)–(2.24) and Hypotheses (H3)–(H5) that

$$\frac{1}{\tau}\left|\int_0^T\int_\Omega (\rho_i^{(\tau)} - \sigma_\tau\rho_i^{(\tau)})\phi_i dxdt\right| \leq C\|\nabla\boldsymbol{v}^{(\tau)}\|_{L^2(\Omega_T)}\|\nabla\phi\|_{L^2(\Omega_T)}$$

$$+ \sum_{i=1}^{n-1} \|M_i^{(\tau)}/\theta^{(\tau)}\|_{L^\infty(\Omega_T)} \|\nabla \log \theta^{(\tau)}\|_{L^2(\Omega_T)} \|\nabla \phi\|_{L^2(\Omega_T)}$$

$$+ \varepsilon \|\boldsymbol{v}^{(\tau)}\|_{L^2(0,T;H^2(\Omega))} \|\phi\|_{L^2(0,T;H^2(\Omega))} + \|\boldsymbol{r}^{(\tau)}\|_{L^2(\Omega_T)} \|\phi\|_{L^2(\Omega_T)}$$

$$\leq C \|\phi\|_{L^2(0,T;H^2(\Omega))},$$

and

$$\frac{1}{\tau} \left| \int_0^T \int_\Omega (E^{(\tau)} - \sigma_\tau E^{(\tau)}) \phi_0 \, dx \, dt \right|$$

$$\leq C + C \|\theta^{(\tau)}\|_{L^{8/3}(\Omega_T)} \|\nabla(\theta^{(\tau)})^2\|_{L^2(\Omega_T)} \|\nabla \phi_0\|_{L^8(\Omega_T)}$$

$$+ \sum_{j=1}^{n-1} \|M_j^{(\tau)}/\theta^{(\tau)}\|_{L^\infty(\Omega_T)} \|\theta^{(\tau)}\|_{L^{8/3}(\Omega_T)} \|\nabla v_j^{(\tau)}\|_{L^2(\Omega_T)} \|\nabla \phi_0\|_{L^8(\Omega_T)}$$

$$+ \lambda \|\theta_0 - \theta^{(\tau)}\|_{L^{8/7}(0,T;L^{8/7}(\partial\Omega))} \|\phi_0\|_{L^8(0,T;L^8(\partial\Omega))}$$

$$+ \varepsilon \|\theta^{(\tau)}\|_{L^3(\Omega_T)} \|\log \theta^{(\tau)}\|_{L^2(0,T;H^2(\Omega))} \|D^2 \phi_0\|_{L^6(\Omega_T)}$$

$$+ \varepsilon \|\theta^{(\tau)}\|_{L^{16/3}(\Omega_T)} \|\nabla \log \theta^{(\tau)}\|_{L^4(\Omega_T)}^3 \|\nabla \phi_0\|_{L^{16}(\Omega_T)}$$

$$+ \varepsilon C \left(1 + \|\theta^{(\tau)} \log \theta^{(\tau)}\|_{L^2(\Omega_T)}\right) \|\phi_0\|_{L^2(\Omega_T)} \leq C \|\phi_0\|_{L^{16}(0,T;W^{2,16}(\Omega))}.$$

Since $|E^{(\tau)} - \sigma_\tau E^{(\tau)}| = \rho^0 |\theta^{(\tau)} - \sigma_\tau \theta^{(\tau)}| \geq \rho_* |\theta^{(\tau)} - \sigma_\tau \theta^{(\tau)}|$, this concludes the proof. $\square$

*Step 5: limit* $(\varepsilon, \tau) \to 0$. Estimates (2.27)–(2.28) allow us to apply the Aubin–Lions lemma in the version of [DJ12]. Thus, there exist subsequences that are not relabeled such that as $(\varepsilon, \tau) \to 0$,

$$\rho_i^{(\tau)} \to \rho_i, \quad \theta^{(\tau)} \to \theta \quad \text{strongly in } L^2(\Omega_T), \ i = 1, \ldots, n-1. \tag{2.29}$$

The $L^\infty(\Omega_T)$ bound for $(\rho_i^{(\tau)})$ and the $L^{16/3}(\Omega_T)$ bound for $(\theta^{(\tau)})$ imply the stronger convergences

$$\rho_i^{(\tau)} \to \rho_i \quad \text{strongly in } L^r(\Omega_T) \text{ for all } r < \infty,$$

$$\theta^{(\tau)} \to \theta \quad \text{strongly in } L^\eta(\Omega_T) \text{ for all } \eta < 16/3.$$

The uniform bounds also imply that, up to subsequences,

$$\rho_i^{(\tau)} \rightharpoonup \rho_i \quad \text{weakly in } L^2(0,T;H^1(\Omega)),$$

$$\theta^{(\tau)} \rightharpoonup \theta \quad \text{weakly in } L^2(0,T;H^1(\Omega)),$$

$$\nabla v_i^{(\tau)} \rightharpoonup \nabla v_i \quad \text{weakly in } L^2(0,T;L^2(\Omega)),$$

$$\tau^{-1}(\rho_i^{(\tau)} - \sigma_\tau \rho_i^{(\tau)}) \rightharpoonup \partial_t \rho_i \quad \text{weakly in } L^2(0,T;H^2(\Omega)'),$$

$$\tau^{-1}(\theta^{(\tau)} - \sigma_\tau \theta^{(\tau)}) \rightharpoonup \partial_t \theta \quad \text{weakly in } L^{16/15}(0,T;W^{2,16}(\Omega)'),$$

where $i = 1, \ldots, n-1$ and $j = 1, \ldots, n$. Moreover, as $(\varepsilon, \tau) \to 0$,

$$\varepsilon \log \theta^{(\tau)} \to 0, \quad \varepsilon v_i^{(\tau)} \to 0 \quad \text{strongly in } L^2(0,T;H^2(\Omega)).$$

At this point, $v_i$ is any limit function; we prove below that $v_i = \log(\rho_i/\rho_n)$.

We deduce from the linearity and boundedness of the trace operator $H^1(\Omega) \hookrightarrow H^{1/2}(\partial\Omega)$ that

$$\theta^{(\tau)} \rightharpoonup \theta \quad \text{weakly in } L^2(0,T;H^{1/2}(\partial\Omega)).$$

Using the compact embedding $H^{1/2}(\partial\Omega) \hookrightarrow L^2(\partial\Omega)$, this gives

$$\theta^{(\tau)} \to \theta \quad \text{strongly in } L^2(0,T;L^2(\partial\Omega)).$$

The a.e. convergence of $\rho_i$ for $i = 1,\ldots,n-1$ implies that, up to a subsequence,

$$\rho_n^{(\tau)} = \rho^0 - \sum_{i=1}^{n-1} \rho_i^{(\tau)} \to \rho^0 - \sum_{i=1}^{n-1} \rho_i =: \rho_n \quad \text{a.e. in } \Omega_T.$$

Next, we prove that $\theta$ and $\rho_i$ are positive a.e. We know already that $\theta^{(\tau)}$ and $\rho_i^{(\tau)}$ are positive in $\Omega_T$. It follows from the $L^\infty(0,T;L^1(\Omega))$ bound for $\log\theta^{(\tau)}$ and the a.e. pointwise convergence $\theta^{(\tau)} \to \theta$ that $\log\theta$ is finite a.e. and therefore $\theta > 0$ a.e. in $\Omega_T$. For the positivity of $\rho_i$, we observe first that there exists a constant $C(n) > 0$ such that for all $z_1,\ldots,z_{n-1} \in \mathbb{R}$,

$$\log\left(1 + \sum_{i=1}^{n-1} e^{z_i}\right) \leq C(n)\left(1 + \sum_{i=1}^{n-1} |z_i|\right).$$

Since $\rho_i^{(\tau)} = \rho^0 \exp(v_i^{(\tau)})/\sum_{j=1}^n \exp(v_j^{(\tau)})$, $\rho^0 \geq \rho_*$, and $v_i^{(\tau)}$ is bounded in $L^1(\Omega)$, this implies for sufficiently small $\delta > 0$ that

$$\text{meas}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\} = \text{meas}\left\{(x,t) : -\log\frac{\rho^0(x)\exp(v_i^{(\tau)}(x,t))}{\sum_{j=1}^n \exp(v_j^{(\tau)}(x,t))} \geq -\log\delta\right\}$$

$$\leq \text{meas}\left\{(x,t) : \sum_{j=1}^n |v_j^{(\tau)}(x,t)| \geq C(1 - \log\delta + \log\rho_*)\right\}$$

$$\leq \frac{C}{-\log\delta}\int_0^T \int_\Omega \sum_{i=1}^n |v_i^{(\tau)}(x,t)|dxdt \leq \frac{C}{-\log\delta}, \quad i = 1,\ldots,n-1.$$

We infer from

$$\text{meas}\left\{\liminf_{(\varepsilon,\tau)\to 0}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\}\right\} \leq \liminf_{(\varepsilon,\tau)\to 0} \text{meas}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\}$$

$$\leq \limsup_{(\varepsilon,\tau)\to 0} \text{meas}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\} \leq \text{meas}\left\{\limsup_{(\varepsilon,\tau)\to 0}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\}\right\},$$

and the pointwise convergence $\rho_i^{(\tau)} \to \rho_i$ that in fact equality holds in the previous chain of inequalities, which means that

$$\text{meas}\{(x,t) : \rho_i(x,t) \leq \delta\} = \lim_{(\varepsilon,\tau)\to 0} \text{meas}\{(x,t) : \rho_i^{(\tau)}(x,t) \leq \delta\} \leq \frac{C}{-\log\delta}$$

and $\rho_i > 0$ a.e. in the limit $\delta \to 0$, where $i = 1, \dots, n-1$. We prove in a similar way for $\rho_n^{(\tau)} = \rho^0/(\sum_{j=1}^{n} \exp(v_i^{(\tau)})) > 0$ that $\rho_n > 0$ a.e.

As $\rho_i^{(\tau)}$ converges a.e. to an a.e. positive limit, we have

$$v_i^{(\tau)} = \log \rho_i^{(\tau)} - \log \rho_n^{(\tau)} \to \log \rho_i - \log \rho_n \quad \text{a.e. in } \Omega_T.$$

Thus $v_i = \log \rho_i - \log \rho_n$. Furthermore, $q_i^{(\tau)} = \log \rho_i^{(\tau)} - \log \theta^{(\tau)} \to \log \rho_i - \log \theta =: q_i$ and

$$(\Pi \boldsymbol{q}^{(\tau)})_i = v_i^{(\tau)} - \frac{1}{n} \sum_{j=1}^{n} v_j^{(\tau)} \to v_i - \frac{1}{n} \sum_{j=1}^{n} v_j =: U_i \quad \text{a.e. in } \Omega_T.$$

This shows that $v_i = q_i - q_n$ and $U_i = (q_i - q_n) - \sum_{j=1}^{n}(q_j - q_n)/n = (\Pi \boldsymbol{q})_i$. The a.e. convergence of $(\Pi \boldsymbol{q}^{(\tau)})$ and the boundedness of $r_i$ by Hypothesis (H5) lead to

$$r_i(\Pi \boldsymbol{q}^{(\tau)}, \theta^{(\tau)}) \to r_i(\Pi \boldsymbol{q}, \theta) \quad \text{strongly in } L^\eta(\Omega_T), \ \eta < \infty.$$

By assumption, $M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})$ and $M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})/\theta^{(\tau)}$ are bounded. Then the strong convergences imply that these sequences are converging in $L^q(\Omega_T)$ for $q < \infty$, and the limits can be identified. Thus,

$$M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \to M_{ij}(\boldsymbol{\rho}, \theta) \quad \text{strongly in } L^q(\Omega_T),$$
$$M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})/\theta^{(\tau)} \to M_j(\boldsymbol{\rho}, \theta)/\theta \quad \text{strongly in } L^q(\Omega_T) \text{ for all } q < \infty.$$

This shows that

$$M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) = \frac{1}{\theta^{(\tau)}} M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})\theta^{(\tau)} \to \frac{1}{\theta} M_j(\boldsymbol{\rho}, \theta)\theta = M_j(\boldsymbol{\rho}, \theta)$$

strongly in $L^\eta(\Omega_T)$ for $\eta < 16/3$. Moreover, taking into account (2.26), we have

$$M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla \frac{1}{\theta^{(\tau)}} = -\frac{M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})}{\theta^{(\tau)}} \nabla \log \theta^{(\tau)} \rightharpoonup \frac{M_j(\boldsymbol{\rho}, \theta)}{\theta} \nabla \log \theta$$

weakly in $L^\eta(\Omega_T)$ for $\eta < 8/3$. Finally, by the weak convergence of $(\nabla \boldsymbol{v}^{(\tau)})$ in $L^2(\Omega_T)$,

$$M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)} \rightharpoonup M_{ij}(\boldsymbol{\rho}, \theta) \nabla v_j \quad \text{weakly in } L^\eta(\Omega_T), \ \eta < 2,$$
$$M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)} \rightharpoonup M_j(\boldsymbol{\rho}, \theta) \nabla v_j \quad \text{weakly in } L^\eta(\Omega_T), \ \eta < 16/11,$$
$$M_j(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla \frac{1}{\theta^{(\tau)}} \rightharpoonup -\frac{1}{\theta^2} M_j(\boldsymbol{\rho}, \theta) \nabla \theta \quad \text{weakly in } L^\eta(\Omega_T), \ \eta < 8/7.$$

These convergences allow us to perform the limit $(\varepsilon, \tau) \to 0$. Finally, we can show as in [Jün15, p. 1980f] that the linear interpolant $\widetilde{\rho}_i^{(\tau)}$ of $\rho_i^{(\tau)}$ and the piecewise constant function $\rho_i^{(\tau)}$ converge to the same limit, which leads to $\rho_i^0 = \widetilde{\rho}_i^{(\tau)}(0) \rightharpoonup \rho_i(0)$ weakly in $H^2(\Omega)'$. Thus, the initial datum $\rho_i(0) = \rho_i^0$ is satisfied in the sense of $H^2(\Omega)'$. Similarly, $(\rho \theta)(0) = \rho^0 \theta^0$ in the sense of $W^{2,16}(\Omega)'$. This finishes the proof.

## 2.5 Proof of Existence – Degenerate Case (Theorem 4)

The proof of Theorem 4 is very similar to that one from Section 2.4, therefore we present only the changes in the proof. Steps 1–3 are the same as in the previous section. Only the estimate of $I_4$ is different:

$$I_4 = \int_\Omega \sum_{i,j=1}^{n-1} M_{ij} \nabla v_i \cdot \nabla v_j dx = \int_\Omega \sum_{i,j=1}^{n} M_{ij} \nabla q_i \cdot \nabla q_j dx \geq \frac{c_M}{n} \int_\Omega \sum_{i=1}^{n} \rho_i |\nabla (\Pi \boldsymbol{q})_i|^2 dx.$$

This gives a uniform estimate for $\int_\Omega \rho_i^{(\tau)} |\nabla (\Pi \boldsymbol{q}^{(\tau)})_i|^2 dx$. We claim that it yields a bound for $\nabla (\rho_i^{(\tau)})^{1/2}$ in $L^2(\Omega_T)$. Indeed, we insert the definitions $q_i^{(\tau)} = \log(\rho_i^{(\tau)}/\theta^{(\tau)})$ and $(\Pi \boldsymbol{q}^{(\tau)})_i = q_i^{(\tau)} - \sum_{j=1}^{n} q_j^{(\tau)}/n = \log \rho_i^{(\tau)} - \sum_{j=1}^{n} (\log \rho_j^{(\tau)})/n$ to find that

$$\sum_{i=1}^{n} \rho_i |\nabla (\Pi \boldsymbol{q}^{(\tau)})_i|^2 = \sum_{i=1}^{n} \rho_i^{(\tau)} \left| \nabla \log \rho_i^{(\tau)} - \frac{1}{n} \sum_{j=1}^{n} \nabla \log \rho_j^{(\tau)} \right|^2$$

$$= \sum_{i=1}^{n} \rho_i^{(\tau)} |\nabla \log \rho_i^{(\tau)}|^2 - \frac{2}{n} \nabla \rho^0 \cdot \sum_{j=1}^{n} \nabla \log \rho_j^{(\tau)} + \frac{\rho^0}{n^2} \left| \sum_{j=1}^{n} \nabla \log \rho_j^{(\tau)} \right|^2$$

$$\geq 4 \sum_{i=1}^{n} |\nabla (\rho_i^{(\tau)})^{1/2}|^2 - 4 |\nabla (\rho^0)^{1/2}|^2.$$

This shows the claim.

In contrast to Step 4 in Section 2.4, we do not have a uniform bound for $v_i^{(\tau)}$ in $L^2(0, T; H^1(\Omega))$ but a bound for $(\rho_i^{(\tau)})^{1/2}$. We deduce from the $L^\infty$ bound for $\rho_i^{(\tau)}$ a bound for $\rho_i^{(\tau)}$ in $L^2(0, T; H^1(\Omega))$, using $\nabla \rho_i^{(\tau)} = (\rho_i^{(\tau)})^{1/2} \nabla (\rho_i^{(\tau)})^{1/2}$. This bound changes the proof of estimate (2.28) for the time translates. In fact, we just have to replace the estimations involving $\nabla v_j^{(\tau)}$:

$$\int_0^T \int_\Omega \left| \sum_{j=1}^{n-1} M_{ij}^{(\tau)} \nabla v_j^{(\tau)} \cdot \nabla \phi_j dx dt \right| dx dt = \int_0^T \int_\Omega \left| \sum_{j=1}^{n} M_{ij}^{(\tau)} \nabla \log \rho_j^{(\tau)} \cdot \nabla \phi_i \right| dx dt$$

$$\leq \sum_{j=1}^{n} \| M_{ij}^{(\tau)} / \rho_j^{(\tau)} \|_{L^\infty(\Omega_T)} \| \nabla \rho_j^{(\tau)} \|_{L^2(\Omega_T)} \| \nabla \phi_i \|_{L^2(\Omega_T)},$$

$$\int_0^T \int_\Omega \left| \sum_{j=1}^{n-1} M_j^{(\tau)} \nabla v_j^{(\tau)} \cdot \nabla \phi_0 \right| dx dt = \int_0^T \int_\Omega \left| \sum_{j=1}^{n} M_j^{(\tau)} \nabla \log \rho_j^{(\tau)} \cdot \nabla \phi_0 \right| dx dt$$

$$\leq \sum_{j=1}^{n} \| M_{ij}^{(\tau)} / \rho_j^{(\tau)} \|_{L^\infty(\Omega_T)} \| \nabla \rho_j^{(\tau)} \|_{L^2(\Omega_T)} \| \nabla \phi_0 \|_{L^2(\Omega_T)}.$$

This yields (2.28).

The $L^2(0, T; H^1(\Omega))$ estimate for $\rho_i^{(\tau)}$ and (2.28) allow us to apply the Aubin–Lions lemma in the version of [DJ12] yielding, up to a subsequence, the strong convergence

$\rho_i^{(\tau)} \to \rho_i$ in $L^2(\Omega_T)$ as $(\varepsilon, \tau) \to 0$ and, because of the boundedness of $\rho_i^{(\tau)}$, in $L^r(\Omega_T)$ for any $r < \infty$.

It remains to perform the limit $(\varepsilon, \tau) \to 0$ in the terms involving $\boldsymbol{v}^{(\tau)}$,

$$\sum_{j=1}^n M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)}, \quad \sum_{i=1}^n M_i(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_i^{(\tau)}, \quad \varepsilon \theta^{(\tau)} (D^2 v_j^{(\tau)} + v_j^{(\tau)}).$$

The last term is easy to treat: The bound for $\sqrt{\varepsilon} v_j^{(\tau)}$ in $L^2(0, T; H^2(\Omega))$ and the strong convergence of $\theta^{(\tau)}$ imply that $\varepsilon \theta^{(\tau)} (D^2 v_j^{(\tau)} + v_j^{(\tau)}) \to 0$ strongly in $L^2(\Omega_T)$. Since $M_{ij}/\rho_j^{(\tau)}$ is bounded by assumption, we have $M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})/\rho_j^{(\tau)} \to M_{ij}(\boldsymbol{\rho}, \theta)/\rho_j$ strongly in $L^r(\Omega_T)$ for $r < \infty$. Hence, using (1.5) and the weak convergence of $(\nabla \rho_j^{(\tau)})$ in $L^2(\Omega_T)$,

$$\sum_{j=1}^{n-1} M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_j^{(\tau)} = \sum_{j=1}^n \frac{M_{ij}(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})}{\rho_j^{(\tau)}} \nabla \rho_j^{(\tau)} \rightharpoonup \sum_{j=1}^n \frac{M_{ij}(\boldsymbol{\rho}, \theta)}{\rho_j} \nabla \rho_j$$

weakly in $L^\eta(\Omega_T)$ for $\eta < 2$. Since $(M_{ij}/\rho_j^{(\tau)}) \nabla \rho_j^{(\tau)}$ is bounded in $L^2(\Omega_T)$, this convergence also holds in $L^2(\Omega_T)$. The limit in the second term $\sum_{i=1}^n M_i(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_i^{(\tau)}$ is performed in an analogous way, leading to

$$\sum_{i=1}^n M_i(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)}) \nabla v_i^{(\tau)} = \sum_{i=1}^n \frac{M_i(\boldsymbol{\rho}^{(\tau)}, \theta^{(\tau)})}{\rho_i^{(\tau)}} \nabla \rho_i^{(\tau)} \rightharpoonup \sum_{i=1}^n \frac{M_i(\boldsymbol{\rho}, \theta)}{\rho_i} \nabla \rho_i$$

weakly in $L^2(\Omega_T)$. This finishes the proof.

# 3 Analysis of a Finite–Volume Scheme for a Quorum Sensing induced Biofilm Model

The results in this chapter are a modified version of the publication [HJZ23].

In this chapter, we provide the discretization and numerical analysis for the quorum sensing induced biofilm model. To this end, we first introduce notation and impose assumptions in section 3.1.1. Afterwards, we provide the finite–volume discretization in section 3.1.2 before presenting the main results of this chapter in section 3.1.3. Having presented the main results, we present the proofs for the existence of discrete solutions in section 3.2, the uniqueness in section 3.3 and the convergence in section 3.5. Finally, we present some numerical simulations in section 3.6, where we show the order of convergence in one dimension and simulations in two dimensions, showing a hollowing effect which was also observed in [ESE17].

While the results in [HJZ23] consider the model of [EPL01], in this thesis we give a modified version of the paper where we perform the computations for the more general model of [EHKE15, ESE17].

## 3.1 Numerical scheme and main results

### 3.1.1 Notation and assumptions

Before we are able to introduce the Finite–Volume Scheme, we need some notation. The following part, which introduces the standard notation for two–point approximation finite–volume methods, can also be found similarly in [DJZ21] and is repeated here for the convenience of the reader.

Let $\Omega \subset \mathbb{R}^2$ be an open, bounded, polygonal domain. We consider only two-dimensional domains, but the generalization to higher space dimensions is straightforward. An admissible mesh of $\Omega$ is given by a family $\mathcal{T}$ of open polygonal control volumes (or cells), a family $\mathcal{E}$ of edges, and a family $\mathcal{P}$ of points $(x_K)_{K \in \mathcal{T}}$ associated to the control volumes and satisfying Definition 9.1 in [EGH00]. This definition implies that the straight line $\overline{x_K x_L}$ between two centers of neighboring cells is orthogonal to the edge $\sigma = K|L$ between two cells. The condition is satisfied, for instance, by triangular meshes whose triangles have angles smaller than $\pi/2$ [EGH00, Example 9.1] or by Voronoï meshes [EGH00, Example 9.2].

The family of edges $\mathcal{E}$ is assumed to consist of interior edges $\mathcal{E}_{\text{int}}$ satisfying $\sigma \subset \Omega$ and boundary edges $\sigma \in \mathcal{E}_{\text{ext}}$ fulfilling $\sigma \subset \partial\Omega$. For a given control volume $K \in \mathcal{T}$, we denote by $\mathcal{E}_K$ the set of edges of $K$. This set splits into $\mathcal{E}_K = \mathcal{E}_{\text{int},K} \cup \mathcal{E}_{\text{ext},K}$. For any $\sigma \in \mathcal{E}$, there exists at least one cell $K_\sigma \in \mathcal{T}$ such that $\sigma \in \mathcal{E}_K$. When $\sigma$ is an interior cell, $\sigma = K|L$, $K_\sigma$ can be either $K$ or $L$.

The admissibility of the mesh and the fact that $\Omega$ is two-dimensional imply that

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathrm{m}(\sigma) \mathrm{d}(x_K, \sigma) \leq 2 \sum_{K \in \mathcal{T}} \mathrm{m}(K) = 2\,\mathrm{m}(\Omega), \tag{3.1}$$

where d is the Euclidean distance in $\mathbb{R}^2$, $\mathrm{m}(\sigma)$ denotes the one-dimensional Lebesgue measure of an edge, and $\mathrm{m}(K)$, $\mathrm{m}(\Omega)$ denote the two-dimensional Lebesgue measure of a cell, the domain, respectively. Let $\sigma \in \mathcal{E}$ be an edge. We define the distance

$$\mathrm{d}_\sigma = \begin{cases} \mathrm{d}(x_K, x_L) & \text{if } \sigma = K|L \in \mathcal{E}_{\mathrm{int},K}, \\ \mathrm{d}(x_K, \sigma) & \text{if } \sigma \in \mathcal{E}_{\mathrm{ext},K}, \end{cases}$$

and introduce the transmissibility coefficient by

$$\tau_\sigma = \frac{\mathrm{m}(\sigma)}{\mathrm{d}_\sigma}. \tag{3.2}$$

We assume that the mesh satisfies the following regularity assumption: There exists $\xi > 0$ such that for all $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}_K$,

$$\mathrm{d}(x_K, \sigma) \geq \xi \mathrm{d}_\sigma. \tag{3.3}$$

The size of the mesh is denoted by $\Delta x = \max_{K \in \mathcal{T}} \mathrm{diam}(K)$.

Let $T > 0$ be the end time, $N_T \in \mathbb{N}$ the number of time steps, $\Delta t = T/N_T$ the time step size, and set $t_k = k\Delta t$ for $k = 0, \ldots, N_T$. We denote by $\mathcal{D}$ an admissible space-time discretization of $\Omega_T := \Omega \times (0,T)$, composed of an admissible mesh $\mathcal{T}$ and the values $(\Delta t, N_T)$. The size of $\mathcal{D}$ is defined by $\chi := \max\{\Delta x, \Delta t\}$.

As it is usual for the finite-volume method, we introduce functions that are piecewise constant in space and time. The finite-volume scheme yields a vector $v_\mathcal{T} = (v_K)_{K \in \mathcal{T}} \in \mathbb{R}^{\#\mathcal{T}}$ of approximate values of a piecewise constant function $v$ such that $v = \sum_{K \in \mathcal{T}} v_K \mathbf{1}_K$, where $\mathbf{1}_K$ is the characteristic function of $K$. We write $v_\mathcal{M} = (v_\mathcal{T}, v_\mathcal{E})$ for the vector that contains the approximate values in the control volumes and on the boundary edges, where $v_\mathcal{E} := (v_\sigma)_{\sigma \in \mathcal{E}_{\mathrm{ext}}} \in \mathbb{R}^{\#\mathcal{E}_{\mathrm{ext}}}$. For such a vector, we use the notation

$$v_{K,\sigma} = \begin{cases} v_L & \text{if } \sigma = K|L \in \mathcal{E}_{\mathrm{int},K}, \\ v_\sigma & \text{if } \sigma \in \mathcal{E}_{\mathrm{ext},K} \end{cases} \tag{3.4}$$

for $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}_K$ and introduce the discrete gradient

$$\mathrm{D}_\sigma v := |\mathrm{D}_{K,\sigma} v|, \quad \text{where } \mathrm{D}_{K,\sigma} v = v_{K,\sigma} - v_K. \tag{3.5}$$

The discrete $H^1(\Omega)$ seminorm and the discrete $H^1(\Omega)$ norm are defined by

$$|v|_{1,2,\mathcal{M}} = \left( \sum_{\sigma \in \mathcal{E}} \tau_\sigma (\mathrm{D}_\sigma v)^2 \right)^{1/2}, \quad \|v\|_{1,2,\mathcal{M}} = \left( \|v\|_{0,2,\mathcal{M}}^2 + |v|_{1,2,\mathcal{M}}^2 \right)^{1/2}, \tag{3.6}$$

where $\|\cdot\|_{0,p,\mathcal{M}}$ denotes the $L^p(\Omega)$ norm

$$\|v\|_{0,p,\mathcal{M}} = \left( \sum_{K \in \mathcal{T}} \mathrm{m}(K) |v_K|^p \right)^{1/p} \quad \text{for } 1 \leq p < \infty.$$

Then, for a given family of vectors $v^k = (v^k_{\mathcal{T}}, v^k_{\mathcal{E}})$ for $k = 1, \ldots, N_T$ and a given nonnegative constant $v^D$ such that $v^k_\sigma = v^D$ for all $\sigma \in \mathcal{E}_{\text{ext}}$, we define the piecewise constant in space and time function $v$ by

$$v(x,t) = \sum_{K \in \mathcal{T}} v^k_K \mathbf{1}_K(x) \quad \text{for } x \in \Omega, \ t \in (t_{k-1}, t_k], \ k = 1, \ldots, N_T. \tag{3.7}$$

For the definition of an approximate gradient for such functions, we need to introduce a dual mesh. Let $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}_K$. The cell $T_{K,\sigma}$ of the dual mesh is defined as follows:

- If $\sigma = K|L \in \mathcal{E}_{\text{int},K}$, then $T_{K,\sigma}$ is that cell ("diamond") whose vertices are given by $x_K$, $x_L$, and the end points of the edge $\sigma$.

- If $\sigma \in \mathcal{E}_{\text{ext},K}$, then $T_{K,\sigma}$ is that cell ("half-diamond") whose vertices are given by $x_K$ and the end points of the edge $\sigma$.

An example of a construction of a dual mesh can be found in [CHLP03]. For an illustration, we refer to [GHHK20]. The cells $T_{K,\sigma}$ define, up to a negligible set, a partition of $\Omega$. The definition of the dual mesh implies the following property. As the straight line between two neighboring centers of cells $\overline{x_K x_L}$ is orthogonal to the edge $\sigma = K|L$, it follows that

$$\text{m}(\sigma)\text{d}(x_K, x_L) = 2\,\text{m}(T_{K,\sigma}) \quad \text{for all } \sigma = K|L \in \mathcal{E}_{\text{int},K}. \tag{3.8}$$

The approximate gradient of a piecewise constant function $v$ in $\Omega_T$ is given by

$$\nabla^{\mathcal{D}} v(x,t) = \frac{\text{m}(\sigma)}{\text{m}(T_{K,\sigma})} \text{D}_{K,\sigma} v^k \nu_{K,\sigma} \quad \text{for } x \in T_{K,\sigma}, \ t \in (t_{k-1}, t_k], , k = 1, \ldots, N_T,$$

where the discrete operator $\text{D}_{K,\sigma}$ is given in (3.5) and $\nu_{K,\sigma}$ is the unit vector that is normal to $\sigma$ and that points outward of $K$.

### 3.1.2 Numerical scheme

We are now in the position to formulate the finite-volume discretization of (1.8)–(1.13). Let $\mathcal{D}$ be an admissible discretization of $\Omega_T$. The initial conditions are discretized by the averages

$$M^0_K = \frac{1}{\text{m}(K)} \int_K M^0(x)dx, \quad N^0_K = \frac{1}{\text{m}(K)} \int_K N^0(x)dx, \tag{3.9}$$

$$S^0_K = \frac{1}{\text{m}(K)} \int_K S^0(x)dx, \quad A^0_K = \frac{1}{\text{m}(K)} \int_K A^0(x)dx. \tag{3.10}$$

for $K \in \mathcal{T}$. On the Dirichlet boundary, we set $M^k_\sigma = M^D$, $N^k_\sigma = 0$, $S^k_\sigma = 1$ and $A^k_\sigma = 0$ for $\sigma \in \mathcal{E}_{\text{ext}}$ at time $t_k$.

Let $M^k_K$, $N^k_K$, $S^k_K$ and $A^k_K$ be some approximations of the mean values of $M(\cdot, t_k)$, $N(\cdot, t_k)$, $S(\cdot, t_k)$ and $A(\cdot, t_k)$, respectively, in the cell $K$. Then the elements $M^k_K$, $N^k_K$, $S^k_K$ and $A^k_K$ are solutions to

$$\frac{\text{m}(K)}{\Delta t}(M^k_K - M^{k-1}_K) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}^k_{M,K,\sigma} = \text{m}(K)g_1(M^k_K, S^k_K, A^k_K), \tag{3.11}$$

$$\frac{m(K)}{\Delta t}(N_K^k - N_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{N,K,\sigma}^k = m(K)g_2(M_K^k, N_K^k, S_K^k, A_K^k), \tag{3.12}$$

$$\frac{m(K)}{\Delta t}(S_K^k - S_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{S,K,\sigma}^k = m(K)g_3(M_K^k, N_K^k, S_K^k), \tag{3.13}$$

$$\frac{m(K)}{\Delta t}(A_K^k - A_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{A,K,\sigma}^k = m(K)g_4(M_K^k, N_K^k, A_K^k), \tag{3.14}$$

the numerical fluxes are defined as

$$\mathcal{F}_{M,K,\sigma}^k = -\tau_\sigma d_1 D_{K,\sigma}F(M^k), \quad \mathcal{F}_{N,K,\sigma}^k = -\tau_\sigma d_2 D_{K,\sigma}N^k, \tag{3.15}$$

$$\mathcal{F}_{S,K,\sigma}^k = -\tau_\sigma d_3 D_{K,\sigma}S^k, \quad \mathcal{F}_{A,K,\sigma}^k = -\tau_\sigma d_4 D_{K,\sigma}A^k, \tag{3.16}$$

where $K \in \mathcal{T}$, $\sigma \in \mathcal{E}_K$, $k \in \{1, \ldots, N_T\}$, and we recall definitions (1.14)–(1.17) for $g_i$, $i = 1, \ldots, 4$, (1.19) for $F$, and (3.2) for $\tau_\sigma$.

For the convenience of the reader, we recall the discrete integration-by-parts formula for piecewise constant functions $v = (v_\mathcal{T}, v_\mathcal{E})$:

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{K,\sigma} v_K = -\sum_{\sigma \in \mathcal{E}} \mathcal{F}_{K,\sigma} D_{K,\sigma}v + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \mathcal{F}_{K,\sigma} v_\sigma, \tag{3.17}$$

where $\mathcal{F}_{K,\sigma}$ is a numerical flux like in (3.15)–(3.16).

### 3.1.3 Main results and Key Ideas

We impose the following hypotheses:

(H1) Domain: $\Omega \subset \mathbb{R}^2$ is a bounded polygonal domain.

(H2) Discretization: $\mathcal{D}$ is an admissible discretization of $\Omega_T := \Omega \times (0, T)$ satisfying the regularity condition (3.3) and $\Delta t < 1/2$.

(H3) Initial data: $S^0$, $M^0 \in L^2(\Omega)$ satisfy $0 \leq S^0 \leq 1$ and $0 \leq M^0 < 1$ in $\Omega$. $N^0, A^0 \in L^\infty$ satisfy $0 \leq N^0$ and $0 \leq A^0 \leq A_{max}$ where $A_{max} > 0$ is a constant such that

$$A_{max} \geq (\alpha + \beta)\frac{1 + 2(\|N^0\|_{L^\infty(\Omega)} + T\eta)\exp(2T)}{\lambda}.$$

(H4) Dirichlet datum: $M^D \in (0, 1)$ is a constant.

(H5) Parameters: $d_i > 0$ for $i = 1, 2, 3, 4$, $k_1, k_2 > 0$, $\alpha, \beta, \eta, \mu > 0$, $n > 1$, $a > 1$, and $b > 0$.

**Remark 12** (Discussion of the hypotheses). Conditions $M^0 < 1$ and $M^D < 1$ allow for the proof of $M_K^k < 1$ for all $K \in \mathcal{T}$ and $k = 1, \ldots, N_T$, thus avoiding quenching of the solution, i.e. the occurrence of regions with $M_K^k = 1$. We assume that $M^D$ is positive to be able to introduce an entropy variable. This condition can be relaxed by introducing

an approximation procedure. The assumption that the boundary biomass is constant is imposed for simplicity. It can be generalized to piecewise constant or time-dependent boundary data, for instance. Moreover, mixed Dirichlet–Neumann boundary conditions or mixed Dirichlet–Robin boundary conditions for the biomass could be imposed as well; see [ESE17, Section 3]. In principle, space-dependent boundary data $M^D$ can be treated with an entropy method for finite-volume schemes as in [CCHGJ19]. The delicate point is here the definition of the entropy variable $W_K^\varepsilon$, which requires a piecewise approximation of the Dirichlet data. We may assume that the diffusion coefficients $d_1, \ldots, d_4$ depend on the spatial variable. In this case, they have to be assumed to be strictly positive. The condition $a > 1$ corresponds to "very fast diffusion". In numerical simulations, usually the values $a = b = 4$ are chosen [EPL01, Table 1]. $\qquad\square$

Our first main result concerns the existence of solutions to the numerical scheme. We introduce the function

$$Z(M) := \int_0^M F(s)ds - F(M^D)(M - M^D), \quad M \in [0,1). \tag{3.18}$$

**Theorem 13** (Existence of discrete solutions). *Assume that Hypotheses (H1)–(H5) hold. Then, for every $k = 1, \ldots, N_T$, there exists a solution $(M^k, N^k, S^k, A^k)$ to scheme (3.9)–(3.16) satisfying*

$$0 \le M_K^k < 1, \quad 0 \le N_K^k \le N_{max}, \quad 0 \le S_K^k \le 1, \quad 0 \le A_K^k \le A_{max} \quad \text{for all } K \in \mathcal{T}, \tag{3.19}$$

*where $N_{max} > 0$ depends only on $N^0$, $T$ and $\eta$ and $A_{max} > 0$ depends only on $\alpha$, $\beta$, $\lambda$ and $N_{max}$. Furthermore there exist positive constants $C_1$ and $C_2$ independent of $\Delta x$ and $\Delta t$ such that*

$$\|Z(M^k)\|_{0,1,\mathcal{M}} + \Delta t C_1 \|F(M^k)\|_{1,2,\mathcal{M}}^2 \le \|Z(M^{k-1})\|_{0,1,\mathcal{M}} + \Delta t C_2. \tag{3.20}$$

*Moreover, if $M^0 \ge m_0$ in $\Omega$ and $M_D \ge m_0$ for some $m_0 > 0$ then any discrete solution $(M^k, N^k, S^k, A^k)$ to scheme (3.9)–(3.16) fulfilling bounds (3.19) satisfies*

$$M_K^k \ge m_0 \exp\left(-\left(k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n}\right) t_k\right) \quad \text{for all } K \in \mathcal{T}, \ k = 1, \ldots, N_T. \tag{3.21}$$

The existence result is proved by a fixed-point argument based on a topological degree result. The main difficulty is to approximate the equations in such a way that the singular point $M = 1$ is avoided. This can be done, as in [EZE09], by introducing a cut-off approximation $f_\varepsilon(M)$ of $f(M)$. Then, by the comparison principle, it is possible to show the bound $M^\varepsilon \le 1 - \delta(\varepsilon)$ for the approximate biomass $M^\varepsilon$, where $\delta(\varepsilon) \in (0,1)$. Since the comparison principle cannot be easily extended to the discrete case, we have chosen another approach.

We first formulate a regularized problem for each time step by introducing the entropy variable $W_K^\varepsilon := Z_\varepsilon'(M_K^\varepsilon)$, where $Z_\varepsilon$ is the sum of $Z(M_K^\varepsilon)$ and $\varepsilon$ times the Boltzmann entropy (see (3.26)), and by adding higher-order terms of $W_K^\varepsilon$. We then solve the regularized problem for $W_K^\varepsilon$ and obtain the biomass fraction by inverting the relation $W_K^\varepsilon = Z_\varepsilon'(M_K^\varepsilon)$. Then $0 < M_K^\varepsilon < 1$ by definition and we can derive a uniform estimate similar to (3.20).

The uniform bound for $F(M^\varepsilon)$ allows us to perform the deregularization limit and to infer that the a.e. limit function $M_K = \lim_{\varepsilon \to 0} M_K^\varepsilon$ satisfies $M_K < 1$ for all $K \in \mathcal{T}$. The positive lower bound for $M^k$ comes from the fact that the source term $g_1(M_K^k, S_K^k, A_K^k)$ is bounded from below by the linear term $-k_2 M_K^k$, and it is proved by a Stampacchia truncation method.

**Theorem 14** (Uniqueness of discrete solutions)**.** *Assume that Hypotheses (H1)–(H5) hold and that there exists a constant $m_0 > 0$ such that $M^0(x) \ge m_0$ for $x \in \Omega$ and $M^D \ge m_0$. Then there exists $\gamma^* > 0$, depending on the data, the mesh, and $m_0$, such that for all $0 < \Delta t < \gamma^*$, there exists a unique solution to scheme* (3.9)–(3.16)*.*

The proof of the theorem is based on a discrete version of the dual method. On the continuous level, the idea is to choose test functions $\psi$ and $\phi$ solving $-\Delta \phi = M_1 - M_2$, $-\Delta \theta = N_1 - N_2$, $-\Delta \psi = S_1 - S_2$ and $-\Delta \zeta = A_1 - A_2$ with homogeneous Dirichlet boundary data, where $(M_1, N_1, S_1, A_1)$ and $(M_2, N_2, S_2, A_2)$ are two solutions to (1.8)–(1.13) with the same initial data, and to exploit the monotonicity of the nonlinearity $F(M)$. On the discrete level, we replace the diffusion equations for $\phi$, $\theta$, $\psi$ and $\zeta$ by the corresponding finite-volume schemes and estimate similarly as in the continuous case. The restriction on the time step size is due to $L^2(\Omega)$ estimates coming from the source terms.

We also prove that our scheme converges to the continuous model, up to a subsequence. For this result, we introduce a family $(\mathcal{D}_m)_{m \in \mathbb{N}}$ of admissible space-time discretizations of $\Omega_T$ indexed by the size $\chi_m = \max\{\Delta x_m, \Delta t_m\}$ of the mesh, satisfying $\chi_m \to 0$ as $m \to \infty$. We denote by $\mathcal{M}_m$ the corresponding meshes of $\Omega$, by $\Delta t_m$ the corresponding time step sizes and we set $\nabla^m := \nabla^{\mathcal{D}_m}$ [JZ22].

**Theorem 15** (Convergence of the scheme)**.** *Assume that the Hypotheses (H1)–(H5) hold. Let $(\mathcal{D}_m)_{m \in \mathbb{N}}$ be a family of admissible meshes satisfying* (3.3) *uniformly and let $(M_m, N_m, S_m, A_m)_{m \in \mathbb{N}}$ be a corresponding sequence of finite-volume solutions to scheme* (3.9)–(3.16) *constructed in Theorem* 13*. Then there exist $(M, N, S, A) \in L^\infty(\Omega_T; \mathbb{R}^2)$, satisfying $F(M) - F(M^D)$, $N$, $S-1$, $A \in L^2(0, T; H_0^1(\Omega))$, and a subsequence of $(M_m, N_m, S_m, A_m)$ (not relabeled) such that, as $m \to \infty$,*

$$M_m \to M, \quad N_m \to N, \quad S_m \to S, \quad A_m \to A \quad \text{a.e. in } \Omega_T,$$
$$\nabla^m F(M_m) \rightharpoonup \nabla F(M), \quad \nabla^m N_m \rightharpoonup \nabla N, \quad \nabla^m S_m \rightharpoonup \nabla S, \quad \nabla^m A_m \rightharpoonup \nabla A \quad \text{weakly in } L^2(\Omega_T).$$

*Moreover, the limit $(M, N, S, A)$ is a weak solution to* (1.8)–(1.13)*, i.e., for all $\psi_1$, $\psi_2$, $\psi_3$, $\psi_4 \in C_0^\infty(\Omega \times [0, T))$,*

$$- \int_0^T \int_\Omega M \partial_t \psi_1 dx dt - \int_\Omega M^0(x) \psi_1(x, 0) dx + d_1 \int_0^T \int_\Omega \nabla F(M) \cdot \nabla \psi_1 dx dt \qquad (3.22)$$
$$= \int_0^T \int_\Omega g_1(M, S, A) \psi_1 dx dt,$$

$$- \int_0^T \int_\Omega N \partial_t \psi_2 dx dt - \int_\Omega N^0(x) \psi_2(x, 0) dx + d_2 \int_0^T \int_\Omega \nabla N \cdot \nabla \psi_2 dx dt \qquad (3.23)$$
$$= \int_0^T \int_\Omega g_2(M, N, S, A) \psi_2 dx dt,$$

$$-\int_0^T \int_\Omega S\partial_t \psi_3 dx dt - \int_\Omega S^0(x)\psi_3(x,0)dx + d_3 \int_0^T \int_\Omega \nabla S \cdot \nabla \psi_3 dx dt \tag{3.24}$$

$$= \int_0^T \int_\Omega g_3(M,N,S)\psi_3 dx dt,$$

$$-\int_0^T \int_\Omega A\partial_t \psi_4 dx dt - \int_\Omega A^0(x)\psi(x,0)dx + d_4 \int_0^T \int_\Omega \nabla A \cdot \nabla \psi_4 dx dt \tag{3.25}$$

$$= \int_0^T \int_\Omega g_4(M,N,A)\psi_4 dx dt,$$

The convergence proof is based on the uniform estimates derived for the proof of Theorem 13 and a discrete compensated compactness technique [ACM17] needed to identify the nonlinear limits. For the limit $m \to \infty$, we use the techniques of [CHLP03]. If uniqueness for the limiting model holds in the class of weak solutions, the whole sequence $(M_m, N_m, S_m, A_m)$ converges. Uniqueness in a smaller class of functions is proved [ESE17, Lemma 3.6], but we have been unable to show the required regularity of the limit $(M, N, S, A)$ from our approximate system, since the time discretization is not compatible with the technique of [ESE17].

## 3.2 Existence of solutions

For the proof of Theorem 13, we proceed by induction. By Hypothesis (H3), $0 \le M_K^0 < 1$, $0 \le N_k^0 \le N_{max}, 0 \le S_K^0 \le 1, 0 \le A_K^0 \le A_{\max}$ holds for $K \in \mathcal{T}$. Let $(M^{k-1}, N^{k-1}, S^{k-1}, A^{k-1})$ satisfy

$$0 \le M_K^{k-1} < 1,$$
$$0 \le N_K^{k-1} \le N_{max},$$
$$0 \le S_K^{k-1} \le 1,$$
$$0 \le A^{k-1} \le A_{max}$$

for all $K \in \mathcal{T}$ and some $k \in \{1, \ldots, N_T\}$. We use the function $Z_\varepsilon : [0,1) \to \mathbb{R}$, defined by

$$Z_\varepsilon(M) = \int_0^M F(s)ds - F(M^D)(M - M^D) + \varepsilon\left(M \log \frac{M}{M^D} + M^D - M\right), \tag{3.26}$$

where $\varepsilon > 0$ and $F(M)$ is given in (1.19).

*Step 1: Definition of a fixed-point problem.* We formulate the problem in terms of the entropy variable and we add a regularization term. Let $R > 0$ and set

$$\mathcal{K}_R := \big\{(W, N, S, A) \in \mathbb{R}^{4\Lambda} : \|W\|_{1,2,\mathcal{M}} < R, \ \|N\|_{0,2,\mathcal{M}} < R, \ \|S\|_{0,2,\mathcal{M}} < R, \ \|A\|_{0,2,\mathcal{M}} < R,$$
$$W_\sigma = 0, \ N_\sigma = 0, S_\sigma = 1, \ A_\sigma = 0 \ \text{ for } \sigma \in \mathcal{E}_{\text{ext}}\big\},$$

where $\Lambda = \#\mathcal{T} + \#\mathcal{E}_{\text{ext}}$. We define the fixed-point mapping $Q : \mathcal{K}_R \to \mathbb{R}^{4\Lambda}$ by $Q(W, N, S, A) = (W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$, where $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ solves for all $K \in \mathcal{T}$,

$$\varepsilon\left(\text{m}(K)W_K^\varepsilon - \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \text{D}_{K,\sigma}W^\varepsilon\right) \tag{3.27}$$

$$= -\frac{\mathrm{m}(K)}{\Delta t}(M_K - M_K^{k-1}) - \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{M,K,\sigma} + \mathrm{m}(K)g_1(M_K, [S_K]_+, [A_K]_+),$$

$$\frac{\mathrm{m}(K)}{\Delta t}(N_K^\varepsilon - N_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{N,K,\sigma} = \mathrm{m}(K)g_2(M_K, [N_K]_+, [S_K]_+, [A_K]_+), \qquad (3.28)$$

$$\frac{\mathrm{m}(K)}{\Delta t}(S_K^\varepsilon - S_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{S,K,\sigma} = \mathrm{m}(K)g_3(M_K, [S_K]_+, [N_K]_+), \qquad (3.29)$$

$$\frac{\mathrm{m}(K)}{\Delta t}(A_K^\varepsilon - A_K^{k-1}) + \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{N,K,\sigma} = \mathrm{m}(K)g_4(M_K, [N_K]_+, [A_K]_+). \qquad (3.30)$$

the fluxes are as in (3.15)–(3.16), $[z]_+ := \max\{0, z\}$, and we impose the Dirichlet boundary conditions $W_\sigma^\varepsilon = 0$, $N_\sigma^\varepsilon = 0$, $S_\sigma^\varepsilon = 1$, $A_\sigma^\varepsilon = 0$ for $\sigma \in \mathcal{E}_{\mathrm{ext}}$. The value $M_K$ is a function of $W_K$, implicitly defined by

$$W_K = Z_\varepsilon'(M_K) = F(M_K) - F(M^D) + \varepsilon \log \frac{M_K}{M^D}, \quad K \in \mathcal{T}. \qquad (3.31)$$

The map $(0, 1) \to \mathbb{R}$, $M_K \mapsto W_K$ is invertible because the function $Z_\varepsilon'$ is increasing. This shows that $M_K$ is well defined and $M_K \in (0, 1)$ for $K \in \mathcal{T}$. The $\varepsilon$-regularization is needed to obtain a bound for $W^\varepsilon$ in the discrete $H^1(\Omega)$ norm. The existence of a unique solution $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ to (3.27)–(3.30) is a consequence of [EGH00, Lemma 9.2].

We claim that $Q$ is continuous. To show this, we first multiply (3.27) by $W_K^\varepsilon$, sum over $K \in \mathcal{T}$, and use the discrete integration-by-parts formula (3.17):

$$\varepsilon \|W^\varepsilon\|_{1,2,\mathcal{M}}^2 = \varepsilon \sum_{K \in \mathcal{T}} \left( \mathrm{m}(K)(W_K^\varepsilon)^2 - \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(W^\varepsilon)W_K^\varepsilon \right) = J_1 + J_2 + J_3, \quad \text{where}$$

$$J_1 = -\sum_{K \in \mathcal{T}} \frac{\mathrm{m}(K)}{\Delta t}(M_K - M_K^{k-1})W_K^\varepsilon,$$

$$J_2 = -\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{M,K,\sigma}W_K^\varepsilon,$$

$$J_3 = \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left( \frac{[S_K]_+}{k_1 + [S_K]_+} - k_2 - \eta\frac{A_K^n}{1 + A_K^n} \right)M_K W_K^\varepsilon.$$

By the Cauchy–Schwarz inequality and the bound $0 < M_K < 1$, we find that

$$|J_1| \leq \frac{2}{\Delta t} \mathrm{m}(\Omega)^{1/2}\|W^\varepsilon\|_{0,2,\mathcal{M}},$$

$$|J_2| \leq \left( \sum_{K \in \mathcal{T}} \frac{1}{\mathrm{m}(K)} \sum_{\sigma \in \mathcal{E}_K} |\mathcal{F}_{M,K,\sigma}|^2 \right)^{1/2}\|W^\varepsilon\|_{0,2,\mathcal{M}},$$

$$|J_3| \leq \left( \frac{1}{k_1 + 1} + k_2 + \eta \right) \mathrm{m}(\Omega)^{1/2}\|W^\varepsilon\|_{0,2,\mathcal{M}}.$$

Because of the assumption $\|W\|_{1,2,\mathcal{M}} < R$, the flux $|\mathcal{F}_{M,K,\sigma}|$ is bounded from above by a constant depending on $R$. This implies that $|J_2| \leq C(R)\|W^\varepsilon\|_{0,2,\mathcal{M}}$, where $C(R) > 0$

is some constant. (Here and in the following, we denote by $C$, $C_i > 0$ generic constants whose value change from line to line.) This shows that $\varepsilon \|W^\varepsilon\|_{1,2,\mathcal{M}} \le C(R)$ for (another) constant $C(R) > 0$. Using similar arguments, we obtain the existence of $C(R) > 0$ such that $\|N^\varepsilon\|_{0,2,\mathcal{M}} \le C(R)$, $\|S^\varepsilon\|_{0,2,\mathcal{M}} \le C(R)$, $\|A^\varepsilon\|_{0,2,\mathcal{M}} \le C(R)$.

Next, let $(W_\ell, N_\ell, S_\ell, A_\ell)_{\ell \in \mathbb{N}} \subset \mathcal{K}_R$ be a sequence satisfying $(W_\ell, N_\ell, S_\ell, A_\ell) \to (W, N, S, A)$ as $\ell \to \infty$. The previous uniform estimates for $(W_\ell^\varepsilon, N_\ell^\varepsilon, S_\ell^\varepsilon, A_\ell^\varepsilon) := Q(W_\ell, N_\ell, S_\ell, A_\ell)$ show that $(W_\ell^\varepsilon, N_\ell^\varepsilon, S_\ell^\varepsilon, A_\ell^\varepsilon)$ is bounded uniformly in $\ell \in \mathbb{N}$. Therefore, there exists a subsequence which is not relabeled such that $(W_\ell^\varepsilon, N_\ell^\varepsilon, S_\ell^\varepsilon, A_\ell^\varepsilon) \to (W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ as $\ell \to \infty$. Taking the limit $\ell \to \infty$ in (3.27)–(3.30), we see that $(W^\varepsilon, S^\varepsilon, N^\varepsilon, A^\varepsilon) = Q(W, S, N, A)$. We deduce from the uniqueness of the limit that the whole sequence converges, which means that $Q$ is continuous.

*Step 2: Brouwer topological degree.* We wish to show that $Q$ admits a fixed point. To this end, we use a topological degree argument [Dei85, Chap. 1] and prove that $\deg(I - Q, \mathcal{K}_R, 0) = 1$, where deg is the Brouwer topological degree. By the properties of the Brouwer topological degree, $\deg(I - Q, \mathcal{K}_R, 0) \ne 0$ implies the existence of $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon) \in \mathcal{K}_R$ such that $(I - Q)(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon) = 0$. Furthermore, by definition of deg, we have $\deg(I, \mathcal{K}_R, 0) = 1$, since $0 \in \mathcal{K}_R$. We deduce from the invariance by homotopy that $\deg(I - \rho Q, \mathcal{K}_R, 0)$ is invariant in $\rho$, if any solution $(S^\varepsilon, W^\varepsilon, \rho) \in \overline{\mathcal{K}}_R \times [0,1]$ to the fixed-point equation $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon) = \rho Q(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ satisfies $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon, \rho) \notin \partial \mathcal{K}_R \times [0,1]$. Therefore, it is sufficient to prove that any solution to the fixed-point equation satisfies $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon, \rho) \notin \partial \mathcal{K}_R \times [0,1]$ for sufficiently large values of $R > 0$.

Let $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon, \rho)$ be a fixed point of $Q$ and assume that $\rho \ne 0$, the case $\rho = 0$ being clear. Then $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ solves

$$\varepsilon \left( \mathrm{m}(K) W_K^\varepsilon - \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma} W^\varepsilon \right) \tag{3.32}$$

$$= -\rho \frac{\mathrm{m}(K)}{\Delta t} (M_K^\varepsilon - M_K^{k-1}) - \rho \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{M,K,\sigma}^\varepsilon + \rho \, \mathrm{m}(K) g_1(M_K^\varepsilon, [S_K^\varepsilon]_+, [A_K^\varepsilon]_+)$$

$$\frac{\mathrm{m}(K)}{\Delta t} (N_K^\varepsilon - \rho N_K^{k-1}) + \rho \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{S,K,\sigma}^\varepsilon = \rho \, \mathrm{m}(K) g_2(M_K^\varepsilon, [N_K^\varepsilon]_+, [S_K^\varepsilon]_+, [A_K^\varepsilon]_+) \tag{3.33}$$

$$\frac{\mathrm{m}(K)}{\Delta t} (S_K^\varepsilon - \rho S_K^{k-1}) + \rho \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{S,K,\sigma}^\varepsilon = \rho \, \mathrm{m}(K) g_3(M_K^\varepsilon, [N_K^\varepsilon]_+, [S_K^\varepsilon]_+), \tag{3.34}$$

$$\frac{\mathrm{m}(K)}{\Delta t} (A_k^\varepsilon - \rho A_K^{k-1}) + \rho \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{A,K,\sigma}^\varepsilon = \rho \, \mathrm{m}(K) g_4(M^\varepsilon, [N^\varepsilon]_+, [A^\varepsilon]_+) \tag{3.35}$$

for $K \in \mathcal{T}$ with the boundary conditions $W_\sigma^\varepsilon = 0$, $N_\sigma^\varepsilon = 0$, $S_\sigma^\varepsilon = 1$, $A_\sigma^\varepsilon = 0$ for $\sigma \in \mathcal{E}_K$, the fluxes are given by (3.15)–(3.16) with $(M, N, S, A)$ replaced by $(M^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$, and $M_K^\varepsilon$ is the unique solution to (3.31) with $W_K$ replaced by $W_K^\varepsilon$.

*Step 3: A priori estimates.* We establish some a priori estimates for the fixed points $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ of $Q$, which are uniform in $R$. Definition (3.31) immediately gives the bound $0 < M_K^\varepsilon < 1$ for all $K \in \mathcal{T}$. Recall that $\lim_{M \nearrow 1} F(M) = +\infty$.

**Lemma 16** (Pointwise bounds for $(N^\varepsilon, S^\varepsilon, A^\varepsilon)$)**.** *There exists $N_{max} > 0$ and $A_{max} > 0$ such that*

$$
\begin{aligned}
0 &\leq N_K^\varepsilon \leq N_{max} \quad \text{for } K \in \mathcal{T}, \\
0 &\leq S_K^\varepsilon \leq 1 \qquad \text{for } K \in \mathcal{T}, \\
0 &\leq A_K^\varepsilon \leq A_{max} \quad \text{for } K \in \mathcal{T}.
\end{aligned}
$$

*Proof.* We start with the lower bounds. To this end, we first multiply (3.33) by $\Delta t[N_K^\varepsilon]_-$, where $[z]_- := \min\{0, z\}$, and sum over $K \in \mathcal{T}$. Using discrete integration by parts, we obtain

$$
\sum_{K \in \mathcal{T}} \mathrm{m}(K)[N_K^\varepsilon]_-^2 + \rho d_2 \Delta t \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma} N^\varepsilon \mathrm{D}_{K,\sigma}[N^\varepsilon]_-
$$
$$
= \rho \sum_{K \in \mathcal{T}} \mathrm{m}(K) N_K^{k-1}[N_K^\varepsilon]_- + \rho \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_2(M_K^\varepsilon, [N_K^\varepsilon]_+, [S_K^\varepsilon]_+, [A_K^\varepsilon]_+)[N_K^\varepsilon]_-.
$$

Using the induction hypothesis, we have $N_K^{k-1} \geq 0$ which gives the nonpositivity of the first term on the right hand side. Furthermore, using $[N_K^\varepsilon]_+[N_K^\varepsilon]_- = 0$ we obtain

$$
g_2(M_K^\varepsilon, [N_K^\varepsilon]_+, [S_K^\varepsilon]_+, [A_K^\varepsilon]_+)[N_K^\varepsilon]_- = \eta \frac{[A_K^\varepsilon]_+^n}{1 + [A_K^\varepsilon]_+^n} M_K^\varepsilon [N_K^\varepsilon]_- \leq 0.
$$

The monotonicity of $z \mapsto [z]_-$ implies

$$
\rho d_2 \Delta t \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma} N^\varepsilon \mathrm{D}_{K,\sigma}[N^\varepsilon]_- \geq 0.
$$

Collecting all estimates, we conclude

$$
\|[N_K^\varepsilon]_-\|_{0,2,\mathcal{M}}^2 \leq 0
$$

and thus $N_K^\varepsilon \geq 0$ for all $K \in \mathcal{T}$. The same argument shows $S_K^\varepsilon \geq 0$ and $A_K^\varepsilon \geq 0$ for all $K \in \mathcal{T}$.

We want to prove the upper bound for $N_K^\varepsilon$.

Let $K_0 \in \mathcal{T}$ with $N_{K_0}^\varepsilon = \max\{N_{K_0}^\varepsilon \mid K \in \mathcal{T}\}$. Then, it is easy to see

$$
\begin{aligned}
\rho \sum_{\sigma \in \mathcal{E}_{K_0}} \mathcal{F}_{N,K_0,\sigma}^\varepsilon &= -\rho \sum_{\sigma \in \mathcal{E}_{K_0}} \tau_\sigma d_2 D_{K_0,\sigma} N_{K_0} \\
&= -\rho \sum_{\sigma \in \mathcal{E}_{K_0}} \tau_\sigma d_2 \left( N_{K_0,\sigma}^\varepsilon - N_{K_0}^\varepsilon \right) \geq 0.
\end{aligned}
$$

Furthermore, we can estimate the right hand side of the equation (3.33) using

$$
\left( \frac{S_K^\varepsilon}{k_1 + S_K^\varepsilon} - k_2 \right) \leq 1
$$

and $0 < \rho \leq 1$. Inserting this in equation (3.33), we find

$$\frac{m(K_0)}{\Delta t} N_{K_0}^\varepsilon \leq \frac{m(K_0)}{\Delta t} N_{K_0}^{k-1} + m(K_0) N_{K_0}^\varepsilon + m(K_0)\eta.$$

Now, dividing by $m(K_0)$ and multiplying with $\Delta t$ we obtain

$$N_{K_0}^\varepsilon \leq N_{K_0}^{k-1} + \Delta t N_{K_0}^\varepsilon + \Delta t \eta \leq N_{K_1}^{k-1} + \Delta t N_{K_0}^\varepsilon + \Delta t \eta,$$

where $N_{K_1}^{k-1} := \max\{N_K^{k-1} \mid K \in \mathcal{T}\}$.

Since $N_{K_0}^\varepsilon$ is the solution at time step $k$, we can iterate the estimate:

$$N_{K_1}^{k-1} \leq N_{K_2}^{k-2} + \Delta t N_{K_1}^{k-1} + \Delta t \eta,$$
$$N_{K_2}^{k-2} \leq N_{K_3}^{k-3} + \Delta t N_{K_2}^{k-2} + \Delta t \eta,$$
$$\vdots$$
$$N_{K_{k-1}}^1 \leq N_{K_k}^0 + \Delta t N_{K_{k-1}}^1 + \Delta t \eta, \tag{3.36}$$

where $N_{K_j}^{k-j} := \max\{N_K^{k-j} \mid K \in \mathcal{T}\}$ for $j = 1, \ldots, k$. We infer

$$N_{K_0}^\varepsilon \leq N_{K_k}^0 + \sum_{\ell=1}^{k} \Delta t N_{K_{k-\ell}}^\ell + (k-1)\Delta t \eta. \tag{3.37}$$

Using $\Delta t < 1/2$, we subtract the term at step $k$ in the inequality (3.37), to arrive at

$$N_{K_0}^\varepsilon \leq 2\left( N_{K_k}^0 + \sum_{\ell=1}^{k-1} \Delta t N_{K_{k-\ell}}^\ell + (k-1)\Delta t \eta \right).$$

Applying Gronwall's inequality and using $(k-1)\Delta t \leq k\Delta t \leq T$, we find

$$N_{K_0}^\varepsilon \leq 2(N_{K_k}^0 + T\eta)\exp(2T) \leq 2(\|N^0\|_{L^\infty(\Omega)} + T\eta)\exp(2T) =: N_{max}.$$

To verify the upper bound for $S^\varepsilon$, we multiply (3.34) by $\Delta t [S_K^\varepsilon - 1]_+$, sum over $K \in \mathcal{T}$, and use discrete integration by parts:

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)\big((S_K^\varepsilon - 1) - (\rho S_K^{k-1} - 1)\big)[S_K^\varepsilon - 1]_+ + \rho d_1 \Delta t \sum_{\sigma \in \mathcal{E}} \mathrm{D}_{K,\sigma}(S^\varepsilon - 1)\mathrm{D}_{K,\sigma}[S^\varepsilon - 1]_+$$
$$= \rho \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_3\big(M_K^\varepsilon, N_K^\varepsilon, S_K^\varepsilon\big)[S_K^\varepsilon - 1]_+ \leq 0, \tag{3.38}$$

since we have always $g_3(M_K^\varepsilon, N_K^\varepsilon, S_K^\varepsilon) \leq 0$. It follows from the induction hypothesis and $\rho \leq 1$ that $\rho S_K^{k-1} \leq 1$, and the first term on the left-hand side can be estimated according to

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)\big((S_K^\varepsilon - 1) - (\rho S_K^{k-1} - 1)\big)[S_K^\varepsilon - 1]_+ \geq \sum_{K \in \mathcal{T}} \mathrm{m}(K)[S_K^\varepsilon - 1]_+^2.$$

We deduce from the monotonicity of $z \mapsto [z]_+$ that the second term on the left-hand side of (3.38) is nonnegative as well. Hence, $\sum_{K \in \mathcal{T}} \mathrm{m}(K)[S_K^\varepsilon - 1]_+^2 \leq 0$ and consequently $S_K^\varepsilon \leq 1$ for all $K \in \mathcal{T}$.

Last but not least, we verify the upper bound for $A$. To this end, we choose $A_{max} > 0$ sufficiently large, such that

$$\lambda A_{max} - [\alpha + \beta](1 + N_{max}) \geq 0.$$

We prove $A_K^\varepsilon \leq A_{max}$ for all $K \in \mathcal{T}$. To this end, we multiply equation (3.35) by $\Delta t[A_K^\varepsilon - A_{max}]_+$ and sum over $K \in \mathcal{T}$, which gives

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K) \left( A_K^\varepsilon - \rho A_K^{k-1} \right) [A_K^\varepsilon - A_{max}]_+ - \Delta t \rho \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathcal{F}_{A,K,\sigma}^\varepsilon [A_K^\varepsilon - A_{max}]_+ \qquad (3.39)$$

$$= \Delta t \rho \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_4(M_K^\varepsilon, N_K^\varepsilon, A_K^\varepsilon)[A_k^\varepsilon - A_{max}]_+. \qquad (3.40)$$

As before, the second term on the left–hand side is nonnegative and

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K) \left( A_K^\varepsilon - \rho A_K^{k-1} \right) [A_K^\varepsilon - A_{max}]_+ \geq \sum_{K \in \mathcal{T}} \mathrm{m}(K)[A_K^\varepsilon - A_{max}]_+^2.$$

It remains to prove the nonpositivity of the right hand side of equation (3.39). If $A_{max} \geq A_K^\varepsilon$, the right–hand side vanishes. If $A_{max} \leq A_K^\varepsilon$, then

$$\lambda A_k^\varepsilon - [\alpha + \beta](1 + N_{max}) \geq \lambda A_{max} - [\alpha + \beta](1 + N_{max}) \geq 0$$

by definition of $A_{max}$. In particular, this implies

$$-\lambda A_K^\varepsilon + [\alpha + \beta](1 + N_{max}) \leq 0.$$

We conclude

$$\begin{aligned}
g_4(M_K^\varepsilon, N_K^\varepsilon, A_K^\varepsilon) &= -\lambda A_K^\varepsilon + \left[ \alpha + \beta \frac{(A_K^\varepsilon)^n}{(1 + (A_K^\varepsilon)^n)} \right] (M_K^\varepsilon + N_K^\varepsilon) \\
&\leq -\lambda A_K^\varepsilon + [\alpha + \beta](1 + N_{max}) \\
&\leq 0,
\end{aligned}$$

and thus

$$\|[A^\varepsilon - A_{max}]_+\|_{0,2,\mathcal{M}}^2 \leq 0,$$

which implies $A_K^\varepsilon \leq A_{max}$ for all $K \in \mathcal{T}$. $\qquad \square$

**Lemma 17** (Estimate for $F(M_K^\varepsilon)$). *Let $\varepsilon \in (0,1]$. Then there exist constants $C_1, C_2 > 0$, depending on $\Omega$ and $M^D$, such that*

$$\varepsilon \Delta t \|W^\varepsilon\|_{1,2,\mathcal{M}}^2 + \rho \|Z(M^\varepsilon)\|_{0,1,\mathcal{M}} + \rho \Delta t C_1 \|F(M^\varepsilon) - F(M^D)\|_{1,2,\mathcal{M}}^2 \qquad (3.41)$$

$$\leq \Delta t C_2 + \|Z_\varepsilon(M^{k-1})\|_{0,1,\mathcal{M}}.$$

*Proof.* We multiply (3.32) by $\Delta t W_K^\varepsilon$, sum over $K$, and use discrete integration by parts with $W_\sigma^\varepsilon = 0$:

$$\varepsilon \Delta t \|W^\varepsilon\|_{1,2,\mathcal{M}}^2 + J_4 + J_5 = J_6, \quad \text{where} \tag{3.42}$$

$$J_4 = \rho \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_K^\varepsilon - M_K^{k-1})W_K^\varepsilon,$$

$$J_5 = \rho \Delta t d_1 \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma} F(M^\varepsilon) \mathrm{D}_{K,\sigma} W^\varepsilon,$$

$$J_6 = \rho \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_1(M_K^\varepsilon, S_K^\varepsilon, A_K^\varepsilon) W^\varepsilon.$$

By the convexity of $Z_\varepsilon$, $(M_K^\varepsilon - M_K^{k-1})Z_\varepsilon'(M_K^\varepsilon) \geq Z_\varepsilon(M_K^\varepsilon) - Z_\varepsilon(M_K^{k-1})$ such that

$$J_4 \geq \rho \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left\{ Z(M_K^\varepsilon) + \varepsilon\left( M_K^\varepsilon \log \frac{M_K^\varepsilon}{M^D} + M^D - M_K^\varepsilon \right) - Z_\varepsilon(M_K^{k-1}) \right\}$$

$$\geq \rho \|Z(M_K^\varepsilon)\|_{0,1,\mathcal{M}} - \rho \|Z_\varepsilon(M_K^{k-1})\|_{0,1\mathcal{M}},$$

where we have used the facts that $Z(M_K^\varepsilon)$ is nonnegative and the function $x \mapsto x \log(x/M^D) + M^D - x$ attains its minimum on $(0,1)$ in $M^D$. The definition of $W_K^\varepsilon$ and the monotonicity of the functions $F$ and log imply that

$$J_5 = \rho \Delta t d_1 \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left([\mathrm{D}_{K,\sigma}(F(M^\varepsilon) - F(M^D))]^2 + \varepsilon \mathrm{D}_{K,\sigma}F(M^\varepsilon)\mathrm{D}_{K,\sigma}\log M^\varepsilon\right) \tag{3.43}$$

$$\geq \rho \Delta t d_1 |F(M^\varepsilon) - F(M^D)|_{1,2,\mathcal{M}}^2 \geq \rho \Delta t d_1 C(\xi)\|F(M^\varepsilon) - F(M^D)\|_{1,2\mathcal{M}}^2,$$

where the last step follows from the discrete Poincaré inequality [BCCHF15, Theorem 3.2]. Recall that by definition of $W$, we have $M_\sigma^\varepsilon = M^D$ since $W_\sigma^\varepsilon = 0$. Finally, by the Young inequality and taking into account the bounds $S_K^\varepsilon \leq 1$ and $M_K^\varepsilon < 1$, we find that

$$J_6 \leq \rho \Delta t \left( \frac{1}{k_1 + 1} + k_2 + \eta \right) \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left( |F(M_K^\varepsilon) - F(M^D)| + \varepsilon M_K^\varepsilon \left| \log \frac{M_K^\varepsilon}{M^D} \right| \right)$$

$$\leq \frac{\delta}{2}\rho \Delta t \left( \frac{1}{k_1 + 1} + k_2 + \eta \right)\|F(M^\varepsilon) - F(M^D)\|_{1,2,\mathcal{M}}^2 + \frac{\Delta t}{2\delta}\left( \frac{1}{k_1 + 1} + k_2 + \eta \right)\mathrm{m}(\Omega)$$

$$+ \varepsilon \Delta t C,$$

where $\delta > 0$, and the constant $C > 0$ may depend on $\Omega$ and $M^D$ but is uniform in $\varepsilon \in (0,1]$. Inserting the estimates for $J_4$, $J_5$, and $J_6$ into (3.42) yields

$$\varepsilon \Delta t \|W^\varepsilon\|_{1,2,\mathcal{M}}^2 + \rho \Delta t \left( d_2 C(\xi) - \frac{\delta}{2}\left( \frac{1}{k_1 + 1} + k_2 + \eta \right) \right)\|F(M^\varepsilon) - F(M^D)\|_{1,2,\mathcal{M}}^2$$

$$+ \rho \|Z(M^\varepsilon)\|_{0,1,\mathcal{M}} \leq \rho \|Z_\varepsilon(M^{k-1})\|_{0,1,\mathcal{M}} + \Delta t C(\delta).$$

Then, choosing $\delta > 0$ sufficiently small shows the conclusion. $\qquad\square$

*Step 4: Concluding the existence of a fixed point.* We deduce from the estimates of Lemmas 16–17 that

$$\|W^\varepsilon\|_{1,2,\mathcal{M}} \leq \frac{1}{\sqrt{\varepsilon \Delta t}}(\|Z_\varepsilon(M^{k-1})\|_{0,1,\mathcal{M}} + \Delta t C)^{1/2}, \quad \|N^\varepsilon\|_{0,2,\mathcal{M}} \leq N_{max}\, \mathrm{m}(\Omega)^{1/2} \quad (3.44)$$

$$\|S^\varepsilon\|_{0,2,\mathcal{M}} \leq \mathrm{m}(\Omega)^{1/2}, \quad \|A^\varepsilon\|_{0,2,\mathcal{M}} \leq A_{max}\, \mathrm{m}(\Omega)^{1/2}. \quad (3.45)$$

Thus, choosing

$$R = \max\left\{ A_{max}\, \mathrm{m}(\Omega)^{1/2}, N_{max}\, \mathrm{m}(\Omega)^{1/2}, \mathrm{m}(\Omega)^{1/2}, \frac{1}{\sqrt{\varepsilon \Delta t}}(\|Z_\varepsilon(M^{k-1})\|_{0,1,\mathcal{M}} + \Delta t C)^{1/2} \right\} + 1,$$

we see that $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon) \notin \partial \mathcal{K}_R$. This implies the invariance in $\rho$, and since we have $\deg(I - \rho Q, \mathcal{K}_R, 0) = 1$ for $\rho = 0$, we infer that $\deg(I - Q, \mathcal{K}_R, 0) = 1$. We conclude that $Q$ admits a fixed point, i.e. a solution $(W^\varepsilon, N^\varepsilon, S^\varepsilon, A^\varepsilon)$ to (3.32)–(3.35).

*Step 5: Limit $\varepsilon \to 0$.* Thanks to Lemmas 16–17 and the bound $0 < M_K^\varepsilon < 1$, there exist subsequences, which are not relabeled, such that $M_K^\varepsilon \to M_K^k$, $N_K^\varepsilon \to N_K^k$, $S_K^\varepsilon \to S_K^k$, $A_K^\varepsilon \to A_K^k$, and $\varepsilon W_K^\varepsilon \to 0$ as $\varepsilon \to 0$ (taking into account (3.44)), where $0 \leq M_K^k \leq 1$, $0 \leq N_K^k \leq N_{max}$, $0 \leq S_K^k \leq 1$ and $0 \leq A_K^k \leq A_{max}$ for all $K \in \mathcal{T}$. Passing to the limit $\varepsilon \to 0$ in (3.41) and taking into account the lower semicontinuity of $F$ (extended for $M = 1$ by setting $F(1) = \infty$), we find that

$$\Delta t C_1 \|F(M^k) - F(M^D)\|_{0,2,\mathcal{M}}^2 \leq \|Z(M^{k-1})\|_{0,1,\mathcal{M}} + \Delta t C < \infty.$$

Thus, $F(M_K^k)$ is finite, which implies that $M_K^k < 1$ for any $K \in \mathcal{T}$. We can perform the limit $\varepsilon \to 0$ in (3.32)–(3.35) to deduce the existence of a solution $(M^k, N^k, S^k, A^k)$ to scheme (3.9)–(3.16).

*Step 6: Positive lower bound for $M^k$.* Again, we proceed by induction. Let $M^0 \geq m_0$ in $\Omega$ and $M^D \geq m_0$. Then $M_K^0 \geq m_0$ for all $K \in \mathcal{T}$. Set

$$m^k := m_0 \left(1 + C\Delta t\right)^{-k},$$

where $C = C(k_2, \eta, A_{max}) := \left( k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right)$. Note that

$$\begin{aligned}
m^k - m^{k-1} &= m_0 \left(1 + C\Delta t\right)^{-k} - m_0 \left(1 + C\Delta t\right)^{-(k-1)} \\
&= (m_0 - m_0(1 + C\Delta t)) \left(1 + C\Delta t\right)^{-k} \\
&= -C\Delta t\, m^k.
\end{aligned}$$

The induction hypothesis reads as $M_K^{k-1} \geq m^{k-1}$ for $K \in \mathcal{T}$. We multiply (3.11) by $\Delta t[M_K^k - m^k]_-$, sum over $K \in \mathcal{T}$, and use discrete integration by parts:

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_K^k - M_K^{k-1})[M_K^k - m^k]_- = J_7 + J_8, \quad \text{where}$$

$$J_7 = -\Delta t \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma} F(M^k) \mathrm{D}_{K,\sigma}[M_K^k - m^k]_-,$$

$$J_8 = \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_1(M^k, S^k, A^k)[M_K^k - m^k]_-.$$

Taking into account that $M_K^{k-1} - m^{k-1} \geq 0$, we estimate the left-hand side according to

$$
\begin{aligned}
\sum_{K \in \mathcal{T}} &\mathrm{m}(K)(M_K^k - M_K^{k-1})[M_K^k - m^k]_- \\
&= \sum_{K \in \mathcal{T}} \mathrm{m}(K)\big((M_K^k - m^k) - (M_K^{k-1} - m^{k-1})\big)[M_K^k - m^k]_- \\
&\quad + \sum_{K \in \mathcal{T}} \mathrm{m}(K)(m^k - m^{k-1})[M_K^k - m^k]_- \\
&\geq \sum_{K \in \mathcal{T}} \mathrm{m}(K)[M_K^k - m^k]_-^2 - C\Delta t m^k \sum_{K \in \mathcal{T}} \mathrm{m}(K)[M_K^k - m^k]_-.
\end{aligned}
$$

Since $F$ and $z \mapsto [z - m^k]_-$ are monotone, we have $J_7 \leq 0$. Furthermore,

$$
\begin{aligned}
J_8 &= \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left(\frac{S_K^k}{k_1 + S_K^k} - k_2 - \eta \frac{(A_K^k)^n}{1 + (A_K^k)^n}\right) M_K^k [M_K^k - m^k]_- \\
&\leq -\left(k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n}\right) \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) M_K^k [M_K^k - m^k]_- \\
&\leq -\left(k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n}\right) \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) m^k [M_K^k - m^k]_-.
\end{aligned}
$$

The terms involving $k_2 + \eta(A_{max})^n/(1 + (A_{max})^n)$ cancel and we end up with

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)[M_K^k - m^k]_-^2 \leq 0.$$

It follows that $[M_K^k - m^k]_- = 0$ and hence $M_K^k \geq m^k \geq m_0 \exp(-Ck\Delta t)$.

## 3.3 Uniqueness of solutions

We proceed by induction. Let $k \in \{1, \ldots, N_T\}$, let $(M_1^k, N_1^k, S_1^k, A_1^k)$ and $(M_2^k, N_2^k, S_2^k, A_2^k)$ be two solutions to scheme (3.9)–(3.16), and assume that $M_1^{k-1} = M_2^{k-1}$, $N_1^{k-1} = N_2^{k-1}$, $S_1^{k-1} = S_2^{k-1}$, $A_1^{k-1} = A_2^{k-1}$. The functions $M_1^k - M_2^k$, $N_1^k - N_2^k$, $S_1^k - S_2^k$ and $A_1^k - A_2^k$ are solutions, respectively, to

$$\frac{\mathrm{m}(K)}{\Delta t}(M_{1,K}^k - M_{2,K}^k) - d_1 \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(F(M_1^k) - F(M_2^k)) = \mathrm{m}(K) G_K^k, \tag{3.46}$$

$$\frac{\mathrm{m}(K)}{\Delta t}(N_{1,K}^k - N_{2,K}^k) - d_2 \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(N_1^k - N_2^k) = \mathrm{m}(K) H_K^k, \tag{3.47}$$

$$\frac{\mathrm{m}(K)}{\Delta t}(S_{1,K}^k - S_{2,K}^k) - d_3 \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(S_1^k - S_2^k) = \mathrm{m}(K) J_K^k, \tag{3.48}$$

$$\frac{\mathrm{m}(K)}{\Delta t}(A_{1,K}^k - A_{2,K}^k) - d_4 \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(A_1^k - A_2^k) = \mathrm{m}(K)L_K^k \tag{3.49}$$

for $K \in \mathcal{T}$, where

$$\begin{aligned}
G_K^k &= \frac{S_{1,K}^k}{k_1 + S_{1,K}} M_{1,K}^k - k_2 M_{1,K}^k - \eta \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n} M_{1,K}^k \\
&\quad - \frac{S_{2,K}^k}{k_1 + S_{2,K}^k} M_{2,K}^k + k_2 M_{2,K}^k + \eta \frac{(A_{2,K}^k)^n}{1 + (A_{2,K}^k)^n} M_{2,K}^k, \\
H_K^k &= \frac{S_{1,K}^k}{k_1 + S_{1,K}^k} N_{1,K}^k - k_2 N_{1,K} + \eta \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n} M_{1,K}^k \\
&\quad - \frac{S_{2,K}^k}{k_1 + S_{2,K}^k} N_{2,K}^k + k_2 N_{2,K} - \eta \frac{(A_{2,K}^k)^n}{1 + (A_{2,K}^k)^n} M_{2,K}^k, \\
J_K^k &= -\frac{\mu S_{1,K}^k}{k_1 + S_{1,K}^k} (M_{1,K}^k + N_{1,K}^k) + \frac{\mu S_{2,K}^k}{k_1 + S_{2,K}^k} (M_{2,K}^k + N_{2,K}^k), \\
L_K^k &= -\lambda A_{1,K}^k + \left[\alpha + \beta \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n}\right] (M_{1,K}^k + N_{1,K}^k) \\
&\quad + \lambda A_{2,K}^k - \left[\alpha + \beta \frac{(A_{2,K}^k)^n}{1 + (A_{2,K}^k)^n}\right] (M_{2,K}^k + N_{2,K}^k).
\end{aligned}$$

Now, let the vectors $(\phi_{\mathcal{T}}^k, \phi_{\mathcal{E}}^k)$, $(\theta_{\mathcal{T}}^k, \theta_{\mathcal{E}}^k)$, $(\psi_{\mathcal{T}}^k, \psi_{\mathcal{E}}^k)$ and $(\zeta_{\mathcal{T}}^k, \zeta_{\mathcal{E}}^k)$ be the unique solutions to

$$\begin{aligned}
-\sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}\phi^k &= \mathrm{m}(K)(M_{1,K}^k - M_{2,K}^k), \\
-\sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}\theta^k &= \mathrm{m}(K)(N_{1,K}^k - N_{2,K}^k) \\
-\sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}\psi^k &= \mathrm{m}(K)(S_{1,K}^k - S_{2,K}^k), \\
-\sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}\zeta^k &= \mathrm{m}(K)(A_{1,K}^k - A_{2,K}^k)
\end{aligned}$$

for $K \in \mathcal{T}$, where we impose the boundary conditions $\phi_\sigma^k = \theta_\sigma^k = \psi_\sigma^k = \zeta_\sigma^k = 0$ for $\sigma \in \mathcal{E}_{\mathrm{ext}}$. The existence and uniqueness of these solutions is a direct consequence of [EGH00, Lemma 9.2]. We multiply (3.46) by $\phi_K^k$ and sum over $K \in \mathcal{T}$:

$$\frac{1}{\Delta t} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_{1,K}^k - M_{2,K}^k)\phi_K^k = I_1 + I_2, \quad \text{where} \tag{3.50}$$

$$I_1 = d_1 \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(F(M_{1,K}^k) - F(M_{2,K}^k))\phi_K^k, \quad I_2 = \sum_{K \in \mathcal{T}} \mathrm{m}(K)G_K^k \phi_K^k.$$

Inserting the equation for $\phi^k$ and using discrete integration by parts gives

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_{1,K}^k - M_{2,K}^k)\phi_K^k = -\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}(\phi^k)\phi_K^k = \sum_{\sigma \in \mathcal{E}} \tau_\sigma (\mathrm{D}_{K,\sigma}\phi^k)^2 = |\phi^k|_{1,2,\mathcal{M}}^2.$$

Concerning the sum $I_1$, we use the equation for $\phi^k$ again, apply discrete integration by parts twice, and take into account the positive lower bound for $M_i^k$ from Theorem 13:

$$\begin{aligned}
I_1 &= d_1 \sum_{K \in \mathcal{T}} (F(M_{1,K}^k) - F(M_{2,K}^k)) \sum_{\sigma \in \mathcal{E}_K} \tau_\sigma \mathrm{D}_{K,\sigma}\phi^k \\
&= -d_1 \sum_{K \in \mathcal{T}} \mathrm{m}(K)(F(M_{1,K}^k) - F(M_{2,K}^k))(M_{1,K}^k - M_{2,K}^k) \\
&\leq -d_1 c_0 \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_{1,K}^k - M_{2,K}^k)^2,
\end{aligned}$$

where $c_0 > 0$ depends on the minimum of $M_1^k$ or $M_2^k$.

To estimate $I_2$, we first rewrite $G_K^k$. We obtain

$$\begin{aligned}
G_K^k &= \left( \frac{S_{1,K}^k}{k_1 + S_{1,K}^k} - k_2 \right)\left( M_{1,K}^k - M_{2,K}^k \right) + \frac{k_1 M_{2,K}^k}{(k_1 + S_{1,K}^k)(k_1 + S_{2,K}^k)}\left( S_{1,K}^k - S_{2,K}^k \right) \\
&\quad - \eta \left[ \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n}\left( M_{1,K}^k - M_{2,K}^k \right) + \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n}M_{2,K}^k \right] + \eta \frac{(A_{2,K}^k)^n}{1 + (A_{2,K}^k)^n}M_{2,K}^k \\
&= \left( \frac{S_{1,K}^k}{k_1 + S_{1,K}^k} - k_2 \right)\left( M_{1,K}^k - M_{2,K}^k \right) + \frac{k_1 M_{2,K}^k}{(k_1 + S_{1,K}^k)(k_1 + S_{2,K}^k)}\left( S_{1,K}^k - S_{2,K}^k \right) \\
&\quad - \eta \left[ \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n}(M_{1,K}^k - M_{2,K}^k) + \frac{(A_{1,K}^k)^n - (A_{2,K}^k)^n}{(1 + (A_{1,K}^k)^n)(1 + (A_{2,K}^k)^n)}M_{2,K}^k \right] \\
&= \left( \frac{S_{1,K}^k}{k_1 + S_{1,K}^k} - k_2 \right)\left( M_{1,K}^k - M_{2,K}^k \right) + \frac{k_1 M_{2,K}^k}{(k_1 + S_{1,K}^k)(k_1 + S_{2,K}^k)}\left( S_{1,K}^k - S_{2,K}^k \right) \\
&\quad - \eta \left[ \frac{(A_{1,K}^k)^n}{1 + (A_{1,K}^k)^n}(M_{1,K}^k - M_{2,K}^k) + \frac{M_{2,K}^k \sum_{\ell=0}^{n-1}(A_{1,K}^k)^\ell (A_{2,K}^k)^{n-1-\ell}}{(1 + (A_{1,K}^k)^n)(1 + (A_{2,K}^k)^n)}(A_{1,K}^k - A_{2,K}^k) \right]
\end{aligned}$$

Finally, because of the bounds $0 \leq M_K^k < 1$, $0 \leq S_K^k \leq 1$ and $0 \leq A_K^k \leq A_{max}$ from Theorem 13, the Young inequality and the discrete Poincaré inequality [BCCHF15, Theorem 3.2],

$$\begin{aligned}
I_2 &\leq -k_2 |\phi^k|_{1,2,\mathcal{M}}^2 + \sum_{K \in \mathcal{T}} \mathrm{m}(K)\left[ \left( \frac{1}{k_1 + 1} + \frac{\eta A_{max}^n}{1 + A_{max}^n} \right) \left| M_{1,K}^k - M_{2,K}^k \right| \right. \\
&\quad \left. + \frac{1}{k_1}\left| S_{1,K}^k - S_{2,K}^k \right| + \eta n A_{max}^{n-1}\left| A_{1,K}^k - A_{2,K}^k \right| \right] \left| \phi_K^k \right| \\
&\leq \frac{3\delta}{4}\left( \frac{1}{k_1 + 1} + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right)^2 \|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \frac{3\delta}{4k_1^2}\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2
\end{aligned}$$

$$+ \frac{3\delta}{4}\eta^2 n^2 A_{max}^{2(n-1)}\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2 + \frac{1}{\delta}\|\phi^k\|_{0,2,\mathcal{M}}^2$$

$$\leq \frac{3\delta}{4}\left(\frac{1}{k_1+1} + \eta\frac{A_{max}^n}{1+A_{max}^n}\right)^2\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \frac{3\delta}{4k_1^2}\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2$$

$$+ \frac{3\delta}{4}\eta^2 n^2 A_{max}^{2(n-1)}\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2 + \frac{C}{\delta\xi}\left|\phi^k\right|_{1,2,\mathcal{M}}^2$$

where $\delta > 0$ is arbitrary. Collecting these estimates, we infer from (3.50) that

$$\left(\frac{1}{\Delta t} - \frac{C}{\delta\xi}\right)|\phi^k|_{1,2,\mathcal{M}}^2 + \frac{3}{4}d_1 c_0\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2$$

$$\leq \frac{3\delta}{4}\left[\left(\frac{1}{k_1+1} + \eta\frac{A_{max}^n}{1+A_{max}^n}\right)^2\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \frac{1}{k_1^2}\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2\right.$$

$$\left. + \eta^2 n^2 A_{max}^{2(n-1)}\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2\right].$$

The same computation gives for equation (3.47)

$$\left(\frac{1}{\Delta t} - \frac{C}{\delta\xi}\right)|\theta^k|_{1,2,\mathcal{M}}^2 + \frac{3}{4}d_2\|N_1^k - N_2^k\|_{0,2,\mathcal{M}}^2$$

$$\leq \frac{3\delta}{4}\left[\left(\frac{1}{k_1+1}\right)^2\|N_1^k - N_2^k\|_{0,2,\mathcal{M}}^2 + \eta^2\left(\frac{A_{max}^n}{1+A_{max}^n}\right)^2\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2\right.$$

$$\left. + \frac{N_{max}^2}{k_1^2}\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2 + \eta^2 n^2 A_{max}^{2(n-1)}\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2\right].$$

Arguing similarly for equation (3.48), we arrive to

$$\left(\frac{1}{\Delta t} - \frac{C}{\delta\xi}\right)|\psi^k|_{1,2,\mathcal{M}}^2 + \frac{3}{4}d_3\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2$$

$$\leq \frac{3\delta}{4}\left[\frac{\mu^2}{(k_1+1)^2}\left(\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \|N_1^k - N_2^k\|_{0,2,\mathcal{M}}^2\right)\right.$$

$$\left. + \frac{\mu^2}{k_1^2}(1 + N_{max})^2\|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2\right].$$

Last but not least, we notice similarly as before that

$$L_k^k = \left(\alpha + \beta\frac{(A_{1,K}^k)^n}{1+(A_{1,K}^k)^n}\right)\left(\left(M_{1,K}^k - M_{2,K}^k\right) + \left(N_{1,K}^k - N_{2,K}^k\right)\right)$$

$$+ \left((M_{2,K}^k + N_{2,K}^k)\frac{\sum_{\ell=0}^{n-1}(A_{1,K}^k)^\ell(A_{2,K}^k)^{n-1-\ell}}{\left(1+(A_{1,K}^k)^n\right)\left(1+(A_{2,K}^k)^n\right)} - \lambda\right)(A_{1,K}^k - A_{2,K}^k).$$

Therefore we can repeat the estimates as before to obtain

$$\left(\frac{1}{\Delta t} - \frac{C}{\delta\xi}\right)|\zeta^k|_{1,2,\mathcal{M}}^2 + \frac{3}{4}d_4\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2$$

$$\leq \frac{3\delta}{4}\left[\left(\alpha + \beta\frac{A_{max}^n}{1+A_{max}^n}\right)^2\left(\|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \|N_1^k + N_2^k\|_{0,2,\mathcal{M}}^2\right)\right.$$

$$\left. + (1+N_{max})^2\left(nA_{max}^{n-1}+\lambda\right)^2\|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2\right]$$

We set

$$R^k := \|M_1^k - M_2^k\|_{0,2,\mathcal{M}}^2 + \|N_1^k - N_2^k\|_{0,2,\mathcal{M}}^2 + \|S_1^k - S_2^k\|_{0,2,\mathcal{M}}^2 + \|A_1^k - A_2^k\|_{0,2,\mathcal{M}}^2$$

and

$$\widetilde{C} := \max\left\{\left(\frac{1}{k_1+1} + \eta\frac{A_{max}^n}{1+A_{max}^n}\right)^2, \frac{1}{k_1^2}, \eta^2 n^2 A_{max}^{2(n-1)}, \frac{\mu^2}{k_1^2}(1+N_{max})^2,\right.$$

$$\left.\left(\alpha + \beta\frac{A_{max}^n}{1+A_{max}^n}\right)^2, (1+N_{max})^2\left(nA_{max}^{n-1}+\lambda\right)^2\right\}$$

Then an addition of the previous inequalities yields

$$\left(\frac{1}{\Delta t} - \frac{C}{\delta\xi}\right)\left(|\phi^k|_{1,2,\mathcal{M}}^2 + |\theta^k|_{1,2,\mathcal{M}}^2 + |\psi^k|_{1,2,\mathcal{M}}^2 + |\zeta^k|_{1,2,\mathcal{M}}^2\right) + \frac{3}{4}\left(\min\{c_0d_1, d_2, d_3, d_4\} - \delta\widetilde{C}\right)R^k \leq 0.$$

Choosing $\delta \leq \widetilde{C}/\min\{c_0d_1, d_2, d_3, d_4\}$ and $\Delta t < C/(\delta\xi)$, both terms are nonnegative, and we infer that $\phi_K^k = \theta_K^k = \psi_K^k = \zeta_K^k = 0$ and consequently

$$M_{1,K}^k - M_{2,K}^k = N_{1,K}^k - N_{2,K}^k = S_{1,K}^k - S_{2,K}^k = A_{1,K}^k - A_{2,K}^k = 0$$

for all $K \in \mathcal{T}$.

## 3.4 Uniform estimates

We establish some estimates for the solution $(M^k, N^k, S^k, A^k)$ constructed in Theorem 13 that are uniform with respect to $\Delta x$ and $\Delta t$. The first bounds follow from the results of Section 3.2.

**Lemma 18** (Uniform estimates I)**.** *There exists a constant $C_3 > 0$ independent of $\Delta x$ and $\Delta t$ such that*

$$0 \leq M_K^k < 1, \quad 0 \leq N_K^k \leq N_{max}, \quad 0 \leq S_K^k \leq 1, \quad 0 \leq A_K^k \leq A_{max} \quad \text{for } K \in \mathcal{T},$$

$$\sum_{k=1}^{N_T}\Delta t\left(\|F(M^k)\|_{1,2,\mathcal{M}}^2 + \|N^k\|_{1,2,\mathcal{M}}^2 + \|S^k\|_{1,2,\mathcal{M}}^2 + \|A^k\|_{1,2,\mathcal{M}}^2\right) \leq C_3.$$

*Proof.* The $L^\infty$ bounds follow directly from Theorem 13, while the discrete gradient bound for $F(M^k)$ is a consequence of Lemma 17. It remains to show the discrete gradient bound

for $N^k$, $S^k$ and $A^k$. We multiply (3.12) by $\Delta t N_K^k$, sum over $K \in \mathcal{T}$, and use discrete integration by parts:

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)(N_K^k - N_K^{k-1})N_K^k = -\Delta t d_2 |N^k|_{1,2,\mathcal{M}}^2 + \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_2(M_K^k, N_K^k, S_K^k, A_K^k)N_K^k, \tag{3.51}$$

where we have used that the boundary terms vanish since $N_\sigma = 0$. The left hand side is bounded from below by

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)(N_K^k - N_K^{k-1})N_K^k \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \mathrm{m}(K) \left( (N_K^k)^2 - (N_K^{k-1})^2 \right),$$

while the second term on the right hand side of equation (3.51) is bounded from above. More precisely, the $L^\infty$–bounds for $(M_K^k, N_K^k, S_K^k, A_K^k)$ imply

$$g_2(M_K^k, N_K^k, S_K^k, A_K^k)N_K^k = \left( \frac{S_K^k}{k_1 + S_K^k} N_K^k - k_2 N_K^k + \eta \frac{(A_K^k)^n}{1 + (A_K^k)^n} \right) N_K^k$$
$$\leq \frac{1}{k_1 + 1} \max\{1, N_{max}^2\} + \eta \frac{A_{max}^n}{1 + A_{max}^n} N_{max}.$$

We conclude

$$\frac{1}{2}\|N^k\|_{0,2,\mathcal{M}}^2 + d_2 \Delta t |N^k|_{1,2,\mathcal{M}}^2 \leq \frac{1}{2}\|N^{k-1}\|_{0,2,\mathcal{M}}^2$$
$$+ \mathrm{m}(\Omega)\Delta t \left( \frac{1}{k_1 + 1}\max\{1, N_{max}^2\} + \eta \frac{A_{max}^n}{1 + A_{max}^n} N_{max} \right),$$

and summation over $k = 1, \ldots, N_T$ gives

$$\frac{1}{2}\|N^{N_T}\|_{0,2,\mathcal{M}}^2 + d_2 \sum_{k=1}^{N_T} \Delta t |N^k|_{1,2,\mathcal{M}}^2 \leq \frac{1}{2}\|N^0\|_{0,2,\mathcal{M}}^2$$
$$+ \mathrm{m}(\Omega)T \left( \frac{1}{k_1 + 1}\max\{1, N_{max}^2\} + \eta \frac{A_{max}^n}{1 + A_{max}^n} N_{max} \right).$$

Due to the different Dirichlet–boundary condition for $S$, we multiply equation (3.13) by $\Delta t(S_K^k - 1)$. Then we proceed as before, i.e. we sum over $K \in \mathcal{T}$ and use discrete integration by parts:

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)(S_K^k - S_K^{k-1})(S_K^k - 1) = -\Delta t d_3 \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma}(S^k) \mathrm{D}_{K,\sigma}(S^k - 1) \tag{3.52}$$
$$+ \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) g_3(M_K^k, N_K^k, S_K^k)(S_K^k - 1)$$
$$\leq -\Delta t d_3 \sum_{\sigma \in \mathcal{E}} \tau_\sigma (\mathrm{D}_{K,\sigma}(S^k - 1))^2 + \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) \frac{\mu S_K^k M_K^k}{k_1 + S_K^k}(M_K^k + N_K^k).$$

Note that the boundary terms vanish since $S_\sigma - 1 = 0$. The left-hand side is bounded from below by

$$\sum_{K \in \mathcal{T}} \mathrm{m}(K)\big((S_K^k - 1) - (S_K^{k-1} - 1)\big)(S_K^k - 1) \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \mathrm{m}(K)\big((S_K^k - 1)^2 - (S_K^{k-1} - 1)^2\big).$$

Using again the upper bounds for $M_K^k$, $N_K^k$ and $S_K^k$, the last term on the right-hand side of (3.52) is bounded by $\Delta t\,\mathrm{m}(\Omega)\mu(1 + N_{max})/(k_1 + 1)$. Therefore, it follows from (3.52) that

$$\frac{1}{2} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(S_K^k - 1)^2 + \Delta t d_3 |S_K^k - 1|_{1,2,\mathcal{M}}^2 \leq \frac{1}{2} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(S_K^{k-1} - 1)^2 + \Delta t\,\mathrm{m}(\Omega)\frac{\mu(1 + N_{max})}{k_1 + 1}.$$

Summing this inequality from $k = 1, \ldots, N_T$, we find that

$$\frac{1}{2}\|S^{N_T} - 1\|_{0,2,\mathcal{M}}^2 + d_3 \sum_{k=1}^{N_T} \Delta t |S_K^k - 1|_{1,2,\mathcal{M}}^2 \leq \frac{1}{2}\|S^0 - 1\|_{0,2,\mathcal{M}}^2 + T\,\mathrm{m}(\Omega)\frac{\mu(1 + N_{max})}{k_1 + 1}.$$

Since $|S_K^k - 1|_{1,2,\mathcal{M}} = |S_K^k|_{1,2,\mathcal{M}}$, this yields the desired estimate. The gradient bound for $A_K^k$ is proved similarly to the gradient bound of $N_K^k$. □

We also need an estimate for the time translates of the solution. For this, let $\phi \in C_0^\infty(\Omega_T)$ be given and define $\phi^k = (\phi_\mathcal{T}^k, \phi_\mathcal{E}^k) \in \mathbb{R}^\Lambda$ (recall that $\Lambda = \#\mathcal{T} + \#\mathcal{E}$) for $k = 1, \ldots, N_T$ by

$$\phi_K^k = \frac{1}{\mathrm{m}(K)} \int_K \phi(x, t_k)dx, \quad \phi_\sigma^k = \frac{1}{\mathrm{m}(\sigma)} \int_\sigma \phi(s, t_k)ds = 0,$$

where $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}_{\text{ext}}$.

**Lemma 19** (Uniform estimates II). *For any $\phi \in C_0^\infty(\Omega_T)$, there exist constants $C_4$, $C_5$, $C_6$, $C_7 > 0$, only depending on the data and the mesh, such that*

$$\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_K^k - M_K^{k-1})\phi_K^k \leq C_4\|\nabla\phi\|_{L^\infty(\Omega_T)},$$

$$\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(N_K^k - N_K^{k-1})\phi_K^k \leq C_5\|\nabla\phi\|_{L^\infty(\Omega_T)},$$

$$\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(S_K^k - S_K^{k-1})\phi_K^k \leq C_6\|\nabla\phi\|_{L^\infty(\Omega_T)},$$

$$\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(A_K^k - A_K^{k-1})\phi_K^k \leq C_7\|\nabla\phi\|_{L^\infty(\Omega_T)}.$$

*Proof.* We only prove the first estimate, since all other estimates are shown similarly. We multiply (3.11) by $\Delta t\phi_K^k$, sum over $K \in \mathcal{T}$ and $k = 1, \ldots, N_T$, and use discrete integration by parts. Then

$$\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_K^k - M_K^{k-1})\phi_K^k = I_3 + I_4, \quad \text{where} \tag{3.53}$$

$$I_3 = -d_1 \sum_{k=1}^{N_T} \Delta t \sum_{\sigma \in \mathcal{E}} \tau_\sigma \mathrm{D}_{K,\sigma} F(M^k) \mathrm{D}_{K,\sigma} \phi^k,$$

$$I_4 = \sum_{k=1}^{N_T} \Delta t \sum_{K \in \mathcal{T}} \mathrm{m}(K) \left( \frac{S_K^k}{k_1 + S_K^k} - k_2 - \eta \frac{(A_K^k)^n}{1 + (A_k^k)^n} \right) M_K^k \phi_K^k.$$

It follows from the Cauchy–Schwarz inequality, Lemma 18, the mean-value theorem, and the mesh regularity (3.3) that

$$|I_3| \le d_1 C \left( \sum_{k=1}^{N_T} \Delta t \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathrm{m}(\sigma) \mathrm{d}_\sigma \left( \frac{D_{K,\sigma} \phi^k}{\mathrm{d}_\sigma} \right)^2 \right)^{1/2}$$

$$\le d_1 C \|\nabla \phi\|_{L^\infty(\Omega_T)} \left( \sum_{k=1}^{N_T} \Delta t \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathrm{m}(\sigma) \mathrm{d}_\sigma \right)^{1/2}$$

$$\le d_1 C \xi^{-1/2} \|\nabla \phi\|_{L^\infty(\Omega_T)} \left( \sum_{k=1}^{N_T} \Delta t \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathrm{m}(\sigma) \mathrm{d}(x_K, \sigma) \right)^{1/2}$$

$$= d_1 C \sqrt{2 \mathrm{m}(\Omega) T \xi^{-1}} \|\nabla \phi\|_{L^\infty(\Omega_T)},$$

where we used (3.1) in the last step. Next, using similar arguments and the discrete Poincaré inequality [BCCHF15, Theorem 3.2],

$$|I_4| \le \left( \frac{1}{k_1 + 1} + k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right) \sqrt{T \mathrm{m}(\Omega)} \left( \sum_{k=1}^{N_T} \Delta t \|\phi^k\|_{0,2,\mathcal{M}}^2 \right)^{1/2}$$

$$\le \left( \frac{1}{k_1 + 1} + k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right) \sqrt{T \mathrm{m}(\Omega) C \xi^{-1}} \left( \sum_{k=1}^{N_T} \Delta t |\phi^k|_{1,2,\mathcal{M}}^2 \right)^{1/2}$$

$$\le \left( \frac{1}{k_1 + 1} + k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right) \sqrt{T \mathrm{m}(\Omega) C \xi^{-1}} \|\nabla \phi\|_{L^\infty(\Omega_T)} \left( \sum_{k=1}^{N_T} \Delta t \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathrm{m}(\sigma) \mathrm{d}_\sigma \right)^{1/2}$$

$$\le C(T, \Omega, \xi) \xi^{-1} \left( \frac{1}{k_1 + 1} + k_2 + \eta \frac{A_{max}^n}{1 + A_{max}^n} \right) \|\nabla \phi\|_{L^\infty(\Omega_T)}.$$

Inserting these estimates into (3.53) shows the first statement of the lemma. □

## 3.5 Convergence of the scheme

The compactness follows from the uniform estimates proved in the previous section and the discrete compensated compactness result obtained in [ACM17, Theorem 3.9].

**Lemma 20** (Compactness). *Let $(M_m, N_m, S_m, A_m)_{m \in \mathbb{N}}$ be a sequence of solutions to scheme (3.9)–(3.16) constructed in Theorem 13. Then there exists $(M, N, S, A) \in L^\infty(\Omega_T; \mathbb{R}^2)$ satisfying $F(M), N, S, A \in L^2(0, T; H^1(\Omega))$ such that, up to a subsequence, as $m \to \infty$,*

$$M_m \to M, \quad N_m \to N, \quad S_m \to S, \quad A_m \to A \quad a.e. \text{ in } \Omega_T,$$

$$F(M_m) \to F(M) \quad \text{strongly in } L^r(\Omega_T) \text{ for } 1 \le r < 2,$$
$$\nabla^m F(M_m) \rightharpoonup \nabla F(M), \quad \nabla^m N_m \rightharpoonup \nabla N, \quad \text{weakly in } L^2(\Omega_T)$$
$$\nabla^m S_m \rightharpoonup \nabla S, \quad \nabla^m A_m \rightharpoonup \nabla A \quad \text{weakly in } L^2(\Omega_T).$$

*Proof.* The a.e. convergence for $M_m$ is a consequence of [ACM17, Theorem 3.9]. Indeed, the estimates in Lemmas 18–19 correspond to conditions (a)–(c) in [ACM17, Prop. 3.8], while assumptions $(A_t 1)$, $(A_x 1)$–$(A_x 3)$ of [ACM17] are satisfied for our implicit Euler finite-volume scheme. In particular, the function $F : [0, 1) \to \mathbb{R}_0^+$ with $F(0) = 0$ and $\lim_{M \nearrow 1} F(M) = +\infty$ is a single-valued maximal monotone graph. Consequently, its inverse $F^{-1}$ is a single-valued maximal monotone graph as well. This allows us to apply [ACM17, Theorem 3.9] so that there exists a subsequence, which is not relabeled, such that $M_m \to M$ and $F(M_m) \to F(M)$ a.e. in $\Omega_T$. In view of Lemma 18, the sequence $(F(M_m))$ is bounded in $L^2(\Omega_T)$. A simple computation shows that $(F(M_m)^r)_{m \in \mathbb{N}}$ is uniformly integrable for $1 \le r < 2$. The a.e. convergence $F(M_m) \to F(M)$ implies the convergence in measure, and thanks to the Vitali's convergence theorem, we conclude that $F(M_m) \to F(M)$ strongly in $L^r(\Omega_T)$ for all $1 \le r < 2$.

As a consequence of the gradient estimate in Lemma 17, there exists a subsequence of $(\nabla^m F(M_m))$ (not relabeled) such that $\nabla^m F(M_m) \rightharpoonup \Psi$ weakly in $L^2(\Omega_T)$ as $m \to \infty$. The limit $\Psi$ can be identified with $F(M)$ by following the arguments in the proof of [CHLP03, Lemma 4.4]. Indeed, the idea is to prove that for all $\phi \in C_0^\infty(\Omega_T; \mathbb{R}^2)$,

$$\mathcal{A}_m := \int_0^T \int_\Omega \nabla^m F(M_m) \cdot \phi \, dx \, dt + \int_0^T \int_\Omega F(M_m) \operatorname{div} \phi \, dx \, dt \to 0$$

as $m \to \infty$. This is done by reformulating the two integrals:

$$\int_\Omega \nabla^m F(M_m) \cdot \phi \, dx = -\frac{1}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int}, K}} \frac{\mathrm{m}(\sigma)}{\mathrm{m}(T_{K,\sigma})} \mathrm{D}_{K,\sigma} F(M_m) \int_{T_{K,\sigma}} \phi(s, t) \cdot \nu_{K,\sigma} \, dx,$$

$$\int_\Omega F(M_m) \operatorname{div} \phi \, dx = \frac{1}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int}, K}} \mathrm{D}_{K,\sigma} F(M_m) \int_\sigma \phi(s, t) \cdot \nu_{K,\sigma} \, ds.$$

Because of the property (see [CHLP03, Lemma 4.4])

$$\left| \frac{1}{\mathrm{m}(T_{K,\sigma})} \int_{T_{K,\sigma}} \phi(t, s) \cdot \nu_{K,\sigma} \, dx - \frac{1}{\mathrm{m}(\sigma)} \int_\sigma \phi(s, t) \cdot \nu_{K,\sigma} \, ds \right| \le \chi_m \|\phi\|_{C^1(\overline{\Omega})}$$

and the uniform estimates for $F(M_m)$ from Lemma 18, it follows that

$$|\mathcal{A}_m| \le \frac{1}{2} \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int}, K}} \mathrm{m}(\sigma) \mathrm{D}_{K,\sigma} F(M^k)$$

$$\times \left| \frac{1}{\mathrm{m}(T_{K,\sigma})} \int_{T_{K,\sigma}} \phi(t, s) \cdot \nu_{K,\sigma} \, dx - \frac{1}{\mathrm{m}(\sigma)} \int_\sigma \phi(s, t) \cdot \nu_{K,\sigma} \, ds \right|$$

$$\le \chi_m C \|\phi\|_{C^1(\overline{\Omega})} \to 0 \quad \text{as } m \to \infty.$$

This implies that $\Psi = \nabla F(M)$. Finally, similar arguments as above show the convergence results for $N_m$, $\nabla^m N_m$, $S_m$, $\nabla^m S_m$ and $A_m$, $\nabla^m A_m$. $\qquad\square$

**Lemma 21** (Convergence of the traces). *Let $(M_m, N_m, S_m, A_m)_{m \in \mathbb{N}}$ be a sequence of solutions to scheme* (3.9)–(3.16) *constructed in Theorem 13. Then the limit function $(M, N, S, A)$ obtained in Lemma 20 satisfies*

$$F(M) - F(M^D),\ N,\ S - 1,\ A \quad \in L^2(0, T; H_0^1(\Omega)).$$

*Proof.* The proof for $N$, $S$ and $A$ is a direct consequence of [BCH13, Prop. 4.9]. For $F(M)$, we follow the proof of [BCH13, Prop. 4.11]. We choose a fixed $m \in \mathbb{N}$ and introduce a definition of the trace of $M_m$, denoted by $\widetilde{M}_m$, such that $\widetilde{M}_m(x, t) = M_K^k$ if $(x, t) \in \sigma \times (t_{k-1}, t_k]$ with $\sigma \in \mathcal{E}_{\text{ext},K}$. Following [BCH13], we wish to prove that

$$\int_0^T \int_{\partial\Omega} (F(\widetilde{M}_m) - F(M))\psi\, dx dt \to 0 \quad \text{as } m \to \infty \tag{3.54}$$

for all $\psi \in C_0^\infty(\partial\Omega \times (0, T))$, from which the claim $F(M) = F(M^D)$ a.e. on $\partial\Omega \times (0, T)$ follows. Indeed, as $M_m = M^D$ on $\partial\Omega \times (0, T)$, we have by the Cauchy–Schwarz inequality,

$$\int_0^T \int_{\partial\Omega} |F(M_m) - F(\widetilde{M}_m)| dx dt = \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} \text{m}(\sigma)|F(M^D) - F(M_K^k)|$$

$$\leq \left( \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} \tau_\sigma |F(M^D) - F(M_K^k)|^2 \right)^{1/2}$$

$$\times \left( \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext},K}} \text{m}(\sigma)\text{d}_\sigma \right)^{1/2}.$$

Hence, thanks to Lemma 18 and the fact that $\text{d}_\sigma = \text{d}(x_K, \sigma) \leq \text{diam}(K) \leq \chi_m$ for every $\sigma \in \mathcal{E}_{\text{ext},K}$, it follows that

$$\int_0^T \int_{\partial\Omega} |F(M_m) - F(\widetilde{M}_m)| dx dt \leq C(T\, \text{m}(\partial\Omega)\chi_m)^{1/2} \to 0 \quad \text{as } m \to \infty,$$

which proves the claim.

Now, as $\Omega$ is assumed to be a polygonal domain, $\partial\Omega$ consists of a finite number of faces denoted by $(\Gamma_i)_{1 \leq i \leq I}$. Similarly to [BCH13, EGHM02], we define for $\delta > 0$ the subset $\Omega_{i,\delta}$ of $\Omega$ such that every $x \in \Omega_{i,\delta}$ satisfies $\text{d}(x, \Gamma_i) < \delta$ and $\text{d}(x, \Gamma_i) < \text{d}(x, \Gamma_j)$ for all $j \neq i$. We also define the subset $\omega_{i,\delta} \subset \Omega_{i,\delta}$ as the largest cylinder of width $\delta$ generated by $\Gamma_i$. Let $\nu_i$ be the unit vector that is normal to $\Gamma_i$, i.e., more precisely, we introduce the set

$$\omega_{i,\delta} := \left\{ x - h\nu_i \in \Omega_i : x \in \Gamma_i,\ 0 < h < \delta \text{ and } [x, x - h\nu_i] \subset \overline{\Omega}_{i,\delta} \right\} \quad \text{for all } 1 \leq i \leq I.$$

Finally, we also introduce the subset $\Gamma_{i,\delta} := \partial\omega_{i,\delta} \cap \Gamma_i$, which fulfills $\text{m}(\Gamma_i \setminus \Gamma_{i,\delta}) \leq C\delta$ for some constant $C > 0$ only depending on $\Omega$.

Let $i \in \{1, \ldots, I\}$ be fixed and let $\psi \in C_0^\infty(\Gamma_i \times (0, T))$. Then there exists $\delta^* = \delta^*(\psi) > 0$ such that for every $\delta \in (0, \delta^*)$, we have $\text{supp}(\psi) \subset \Gamma_{i,\delta} \times (0, T)$. We write

$$\int_0^T \int_{\Gamma_i} (F(\widetilde{M}_m) - F(M))\psi\, dx dt = B_{1,m,\delta} + B_{2,m,\delta} + B_{3,\delta}, \quad \text{where}$$

$$B_{1,m,\delta} = \int_0^T \frac{1}{\delta} \int_{\Gamma_{i,\delta}} \int_0^\delta \left( F(\widetilde{M}_m(x,t)) - F(M_m(x - h\nu_i, t)) \right) \psi(x,t) dh dx dt,$$

$$B_{2,m,\delta} = \int_0^T \frac{1}{\delta} \int_{\Gamma_{i,\delta}} \int_0^\delta \left( F(M_m(x - h\nu_i, t)) - F(M(x - h\nu_i, t)) \right) \psi(x,t) dh dx dt,$$

$$B_{3,\delta} = \int_0^T \frac{1}{\delta} \int_{\Gamma_{i,\delta}} \int_0^\delta \left( F(M(x - h\nu_i, t)) - F(M) \right) \psi(x,t) dh dx dt.$$

We apply the Cauchy–Schwarz inequality to the first term and then use [BCH13, Lemma 4.8] and Lemma 18 to find that

$$|B_{1,m,\delta}| \leq \left( \int_0^T \frac{1}{\delta} \int_{\Gamma_{i,\delta}} \int_0^\delta \left( F(\widetilde{M}_m(x,t)) - F(M_m(x - h\nu_i, t)) \right)^2 dh dx dt \right)^{1/2}$$

$$\times \left( \int_0^T \int_{\Gamma_i} \psi(x,t)^2 dx dt \right)^{1/2} \leq \sqrt{\delta + \chi_m} \|F(M_m)\|_{1,2,\mathcal{M}} \|\psi\|_{L^2(\Gamma_i \times (0,T))}.$$

Taking into account that Lemma 20 implies that $F(M_m) \to F(M)$ strongly in $L^r(\Omega_T)$ for $1 \leq r < 2$, we infer that the second term $B_{2,m,\delta}$ converges to zero as $m \to \infty$. This shows that

$$\lim_{m\to\infty} \left| \int_0^T \int_{\Gamma_i} (F(\widetilde{M}_m) - F(M)) \psi dx dt \right| \leq C\sqrt{\delta} + |B_{3,\delta}|.$$

Since $F(M) \in L^2(0, T; H^1(\Omega))$, the function $F(M)$ has a trace in $L^2(\partial\Omega \times (0, T))$ so that $B_{3,\delta} \to 0$ as $\delta \to 0$. Hence, performing the limit $\delta \to 0$, we conclude that (3.54) holds, finishing the proof. $\square$

It remains to verify that the limit function $(M, N, S, A)$ obtained in Lemma 20 is a weak solution to (1.10)–(1.18). We follow the ideas of [CHLP03] and prove that $M$ solves (3.22), as the proof of (3.24) is analogous. Let $\phi \in C_0^\infty(\Omega \times [0, T))$ and let $\chi_m = \max\{\Delta x_m, \Delta t_m\}$ be sufficiently small such that $\text{supp}(\phi) \subset \{x \in \Omega : \text{d}(x, \partial\Omega) > \chi_m\} \times (0, T)$. The aim is to prove that $F_{10}^m + F_{20}^m + F_{30}^m \to 0$ as $m \to \infty$, where

$$F_{10}^m = -\int_0^T \int_\Omega M_m \partial_t \phi dx dt - \int_\Omega M_m(x, 0) \phi(x, 0) dx,$$

$$F_{20}^m = d_1 \int_0^T \int_\Omega \nabla^m F(M_m) \cdot \nabla \phi dx dt,$$

$$F_{30}^m = -\int_0^T \int_\Omega g_1(M_m, S_m, A_m) \phi dx dt.$$

The convergence results from Lemma 20 allow us to perform the limit $m \to \infty$ in these integrals, leading to

$$F_{10}^m + F_{20}^m + F_{30}^m \to -\int_0^T \int_\Omega M \partial_t \phi dx dt - \int_\Omega M^0(x) \phi(x, 0) dx$$

$$+ d_1 \int_0^T \int_\Omega \nabla F(M) \cdot \nabla \phi dx dt - \int_0^T \int_\Omega g_1(M, S, A) \phi dx dt.$$

Now we set $\phi_K^k = \phi(x_K, t_k)$, multiply (3.11) by $\Delta t \phi_K^{k-1}$, and sum over $K \in \mathcal{T}$ and $k = 1, \ldots, N_T$:

$$F_1^m + F_2^m + F_3^m = 0, \quad \text{where} \tag{3.55}$$

$$F_1^m = \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)(M_K^k - M_K^{k-1})\phi_K^{k-1},$$

$$F_2^m = -d_1 \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int},K}} \tau_\sigma \mathrm{D}_{K,\sigma} F(M^k)\phi_K^{k-1},$$

$$F_3^m = -\sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \mathrm{m}(K)g_1(M_K^k, S_K^k, A_K^k)\phi_K^{k-1}.$$

We claim that $F_{j0}^m - F_j^m \to 0$ as $m \to \infty$ for $j = 1, 2, 3$. Then (3.55) implies that $F_{10}^m + F_{20}^m + F_{30}^m \to 0$ for $m \to \infty$, finishing the proof.

For the first limit, we argue as in [CHLP03, Theorem 5.2]:

$$F_{10}^m = -\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)M_{m,K}^k(\phi_K^k - \phi_K^{k-1}) - \sum_{K \in \mathcal{T}} \mathrm{m}(K)M_{m,K}^0 \phi_K^0$$

$$= -\sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \int_{t_{k-1}}^{t_k} \int_K M_{m,K}^k \partial_t \phi(x_K, t) dx dt - \sum_{K \in \mathcal{T}} \int_K M_{m,K}^0 \phi(x_K, 0) dx.$$

This shows that $|F_{10}^m - F_1^m| \le C\|\phi\|_{C^2(\overline{\Omega_T})} \chi_m \to 0$ as $m \to \infty$.

Next, we use discrete integration by parts to rewrite $F_2^m$:

$$F_2^m = d_1 \sum_{k=1}^{N_T} \Delta t_m \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int},K}} \tau_\sigma \mathrm{D}_{K,\sigma} F(M^k) \mathrm{D}_{K,\sigma} \phi^{k-1}.$$

By the definition of the discrete gradient, we can also rewrite $F_{20}^m$:

$$F_{20}^m = d_1 \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{int,K}} \mathrm{D}_{K,\sigma} F(M^k) \frac{\mathrm{m}(\sigma)}{\mathrm{m}(T_{K,\sigma})} \int_{t_{k-1}}^{t_k} \int_{T_{K,\sigma}} \nabla \phi \cdot \nu_{K,\sigma} dx dt.$$

Hence, using [CHLP03, Theorem 5.1] and the Cauchy–Schwarz inequality, we find that

$$|F_{20}^m - F_2^m| \le d_1 \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int},K}} \mathrm{m}(\sigma) \mathrm{D}_\sigma F(M^k)$$

$$\times \left| \int_{t_{k-1}}^{t_k} \left( \frac{1}{\mathrm{m}(T_{K,\sigma})} \int_{T_{K,\sigma}} \nabla \phi \cdot \nu_{K,\sigma} dx - \frac{1}{\mathrm{d}_\sigma} \mathrm{D}_{K,\sigma} \phi^{k-1} dx \right) dt \right|$$

$$\le d_1 \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\mathrm{int},K}} \mathrm{m}(\sigma) \mathrm{D}_\sigma F(M^k) \cdot C \Delta t_m \chi_m$$

$$\leq C\chi_m d_1 \Big(\sum_{k=1}^{N_T} \Delta t_m \sum_{\sigma \in \mathcal{E}} \mathrm{m}(\sigma)\mathrm{d}_\sigma\Big)^{1/2} \Big(\sum_{k=1}^{N_T} \Delta t_m |F(M^k)|_{1,2,\mathcal{M}}^2\Big)^{1/2}$$

$$\leq C\chi_m d_1 \xi^{-1/2} \Big(\sum_{k=1}^{N_T} \Delta t_m \sum_{\sigma \in \mathcal{E}} \mathrm{m}(\sigma)\mathrm{d}(x_K,\sigma)\Big)^{1/2},$$

where we used the mesh regularity (3.3) in the last step. Taking into account the estimate for $F(M_m)$ from Lemma 17 and the property (3.1), we infer that $F_{20}^m - F_2^m \to 0$.

Finally, using the regularity of $\phi$, we obtain

$$|F_{30}^m - F_3^m| \leq \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K)|g_1(M_K^k, S_K^k, A_K^k)| \left| \int_{t_{k-1}}^{t_k} \left(\phi_K^{k-1} - \frac{1}{\mathrm{m}(K)}\int_K \phi dx\right)dt \right|$$

$$\leq \left(\frac{1}{k_1+1} + k_2 + \eta\frac{A_{max}^n}{1+A_{max}^n}\right) \sum_{k=1}^{N_T} \sum_{K \in \mathcal{T}} \mathrm{m}(K) \left| \int_{t_{k-1}}^{t_k} \left(\phi_K^{k-1} - \frac{1}{\mathrm{m}(K)}\int_K \phi dx\right)dt \right|$$

$$\leq \left(\frac{1}{k_1+1} + k_2 + \eta\frac{A_{max}^n}{1+A_{max}^n}\right) \mathrm{m}(\Omega)T\|\nabla\phi\|_{L^\infty(\Omega_T)}\chi_m \to 0.$$

This finishes the proof.

## 3.6 Numerical experiments

We present in this section some numerical experiments for the biofilm model (3.9)–(3.16) in one and two space dimensions.

### 3.6.1 Implementation of the scheme

The finite-volume scheme (3.9)–(3.16) is implemented in MATLAB. As we used the same method to implement the scheme as [JZ21], we repeat the description of the implementation and the adaptive time step procedure for the convenience of the reader:

Since the numerical scheme is implicit in time, we have to solve a nonlinear system of equations at each time step. In the one-dimensional case, we use Newton's method.

Starting from $(M^{k-1}, N^{k-1}, S^{k-1}, A^{k-1})$, we apply a Newton method with precision $\varepsilon = 10^{-10}$ to approximate the solution to the scheme at time step $k$. In the two-dimensional case, we use a Newton method complemented by an adaptive time-stepping strategy to approximate the solution of the scheme at time $t_k$. More precisely, starting again from $(M^{k-1}, N^{k-1}, S^{k-1}, A^{k-1})$, we launch a Newton method. If the method does not converge with precision $\varepsilon = 10^{-8}$ after at most 50 steps, we multiply the time step by a factor 0.2 and restart the Newton method. At the beginning of each time step, we increase the value of the previous time step size by multiplying it by 1.1. Moreover, we impose the condition $10^{-8} \leq \Delta t_k \leq 10^{-2}$ with an initial time step size equal to $10^{-5}$. Our adaptive time-step strategy aims to improve the numerical performance of our scheme in terms of number of time steps, CPU time, etc. However, this strategy is not mandatory and, as in our one-dimensional test case, we can always implement our scheme with a constant time step with a reasonable size.

### 3.6.2 Test case 1: Rate of convergence in space

We illustrate the order of convergence in space and at final time $T = 10^{-3}$ for the biofilm model in one space dimension with $\Omega = (0, 1)$.

| Parameter | $d_2$ | $d_3$ | $d_4$ | $k_1$ | $k_2$ | $\alpha$ | $\beta$ | $\lambda$ | $\eta$ | $\nu$ | $n$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Value | 4.1667 | 4.1667 | 3.234 | 0.4 | 0.067 | 30.7 | $10\alpha$ | 0.02218 | 0.6 | 793.65 | 2.5 |

Table 3.1: Parameters used in the numerical simulations

Considering remark 1, we assume only homogeneous Dirichlet boundary data for the biomass in our numerics, i.e. $M^D = 0$. For $d_1$, we choose $d_1 = 4.1667$. Except for $d_1$, these values are the same as those considered in [ESE17]. Indeed, in [ESE17] the value of $d_1$ is set to $4.2 \cdot 10^{-8}$. However, with this rather small value of $d_2$, we have to compute the solution of (3.9)–(3.16) for a large final time to obtain a relevant approximation of the order of convergence in space at $T$. Here, the idea was to speed up the dynamics of $M$ to reduce the computational time. We take $a = 2$ and $b = 1$ such that,

$$F(x) = \log(1 - x) + \frac{1}{1 - x} - 1,$$
$$F'(x) = \frac{x}{(1 - x)^2}.$$

Finally, we impose the initial data $N^0(x) = \sin(\pi x)$, $S^0(x) = 1 - 0.2\sin(\pi x)$, $A^0(x) = \sin(\pi x)$ and

$$M^0(x) = 0.2\, g(x - 0.38) + 0.9\, g(x - 0.62),$$
$$\text{where } g(x) = \max\{1 - 9^2 x^2, 0\}.$$

Since exact solutions to the biofilm model are not explicitly known, we compute a reference solution $(M_{\text{ref}}, N_{\text{ref}}, S_{\text{ref}}, A_{\text{ref}})$ on a uniform mesh composed of 5120 cells and with $\Delta t = (1/2\,560)^2 \approx 1.53 \cdot 10^{-7}$. We use this rather small value of $\Delta t$ because the Euler discretization in time exhibits a first-order convergence rate, while we expect a second-order convergence rate in space for scheme (3.9)-(3.16), due to the two-point flux approximation scheme used in this work [DJZ21]. We compute approximate solutions on uniform meshes made of 40, 80, 160, 320, 640, 1280 and 2560 cells, respectively. In Figure 3.1, we present the $L^1(\Omega)$ norm of the difference between the approximate solutions and the average of the reference solution $(M_{\text{ref}}, N_{\text{ref}}, S_{\text{ref}}, A_{\text{ref}})$ at the final time $T = 10^{-3}$. We observe a convergence of order 2 (approximately) for $M$, $N$, $S$ and $A$, respectively, in the $L^1$ norm.

In the $L^2$ norm, we observe a convergence of order $\approx 1.7$ for $M$ and a second order convergence for $N$, $S$ and $A$.
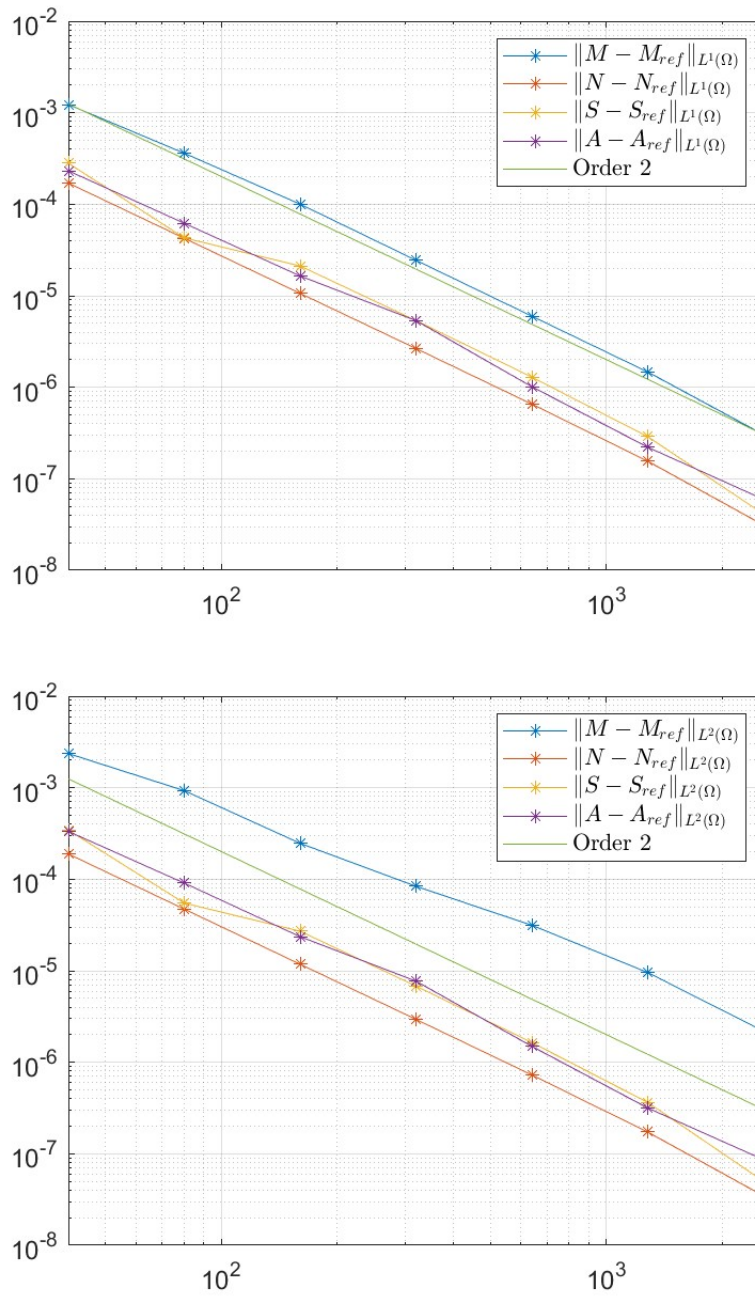
Figure 3.1: Test case 1: $L^1$ and $L^2$ norm of the error between the reference solution and the solutions computed on coarser grids at final time $T = 10^{-3}$.

### 3.6.3 Test case 2: Microbial floc I

We consider the domain $\Omega = (0,1) \times (0,1)$ and the microbial floc which we have considered in [HJZ23]. Therefore, we choose $a = b = 4$ such that

$$F(x) = -\frac{18x^2 - 30x + 13}{3(x-1)^3} + x + 4\log(1-x) - \frac{13}{3}$$

and the initial data $N^0(x,y) = 0$, $S^0(x,y) = 1$, $A^0(x,y) = 0$ and

$$M^0(x,y) = 0.3\,p(x - 0.4, y - 0.5) + 0.9\,p(x - 0.6, y - 0.5),$$
$$\text{where } p(x,y) = \max\{1 - 8^2 x^2 - 8^2 y^2, 0\}.$$

As in test case 1 we take the values from table 3.1, i.e. the values from [ESE17]. However, in the 2D test cases we also choose $d_1$ as in [ESE17], such that $d_1 = 4.2 \cdot 10^{-8}$.

In Figures 3.2–3.5, we illustrate the behavior $M$,$N$,$S$ and $A$ along time for a mesh of $\Omega = (0,1)^2$ composed of 896 triangles. We observe, as in [EPL01,EZE09,ESE17], that after a transient time, the two colonies merge. Due to the high concentration of biomass after this stage, we observe an expansion of the biomass due to the porous-medium type degeneracy. Then, at $T = 10$ we can see in figure 3.5, that the biofilm produced a comparatively large amount of the quorum sensing signal molecule $A$. As described in [ESE17], the quorum sensing signal molecule can switch the biofilm development and put the biofilm growth to a hold. At the same time, the overall number of dispersed cells $N$ is continuously growing due to the increase of signal molecules. As time progresses, we can see that the interplay between the held growth of biofilm and the deteachment of cells causes a hollowing effect similarly to the observations in [ESE17]. Furthermore, we can see in figure 3.2 at time $T = 15$, that the growth now takes place only in areas, where the concentration of the signal molecule (see figure 3.5) is low. At the final time which we consider, namely $T = 25$ we observe in figures 3.3 and 3.5, that the signal molecule as well as the concentration of dispersed cells decreased. Thus, the biomass started to grow again in its new boundaries. Due to the hollowing effect, the growth in the former "center" of the biomass is still put to a hold.

(a) $M$ at $T = 0.0002$

(b) $M$ at $T = 5$

(c) $M$ at $T = 10$

(d) $M$ at $T = 15$

(e) $M$ at $T = 20$

(f) $M$ at $T = 25$

Figure 3.2: Test case 2: Evolution of $M$.

(a) $N$ at $T = 0.0002$

(b) $N$ at $T = 5$

(c) $N$ at $T = 10$

(d) $N$ at $T = 15$

(e) $N$ at $T = 20$

(f) $N$ at $T = 25$

Figure 3.3: Test case 2: Evolution of $N$.

(a) $S$ at $T = 0.0002$

(b) $S$ at $T = 5$

(c) $S$ at $T = 10$

(d) $S$ at $T = 15$

(e) $S$ at $T = 20$

(f) $S$ at $T = 25$

Figure 3.4: Test case 2: Evolution of $S$.

(a) $A$ at $T = 0.0002$

(b) $A$ at $T = 5$

(c) $A$ at $T = 10$

(d) $A$ at $T = 15$

(e) $A$ at $T = 20$

(f) $A$ at $T = 25$

Figure 3.5: Test case 2: Evolution of $A$.

### 3.6.4 Test case 3: Microbial floc II

As mentioned before, we take the values of table 3.1 and $d_1 = 4.2 \cdot 10^{-8}$ as in [ESE17]. As in the second test case, we take $a = b = 4$ and $N^0(x,y) = 0$, $S^0(x,y) = 1$, $A^0(x,y) = 0$. However, as the initial data for the biomass, we consider similarly to [ESE17, 4.1 Microbial floc] a microbial floc which is $M^0 = 0.1$ in a small circular region of the domain. More precisely, we choose

$$M^0(x,y) = \begin{cases} 0.1, & \text{if } \sqrt{(x - 1/2)^2 + (y - 1/2)^2} \leq 1/16, \\ 0, & \text{otherwise.} \end{cases}$$

Note,that in [ESE17] the circular floc takes 0.03% of the domain $[0,1] \times [0,1]$. This is not the case for us, since by definition of the initial data $M^0$, we occupy approximately 0.012% of the domain. We still expect to observe a similar behavior as in [ESE17]. As in case 2, we consider a mesh consisting of 896 triangles. The biomass stays (apart from growth) for $T = 5$ and $T = 10$. At $T = 15$ however, we can already observe that the signal molecule concentration has increased significantly (Figure 3.9) and the detachment of cells has started in the center of the biomass effect (Figure 3.6). Unlike in [ESE17], the biomass does not increase in the inside of the hollowing effect at time $T = 25$, which is probably caused by the different choice of the initial data. Due to the different initial data, the maximum of the signal molecule seems to be reached only at a later time. In this way, unlike in [ESE17], we observe at $T = 25$ only a growth at the edges of the biomass region, but not yet at the void in the center.
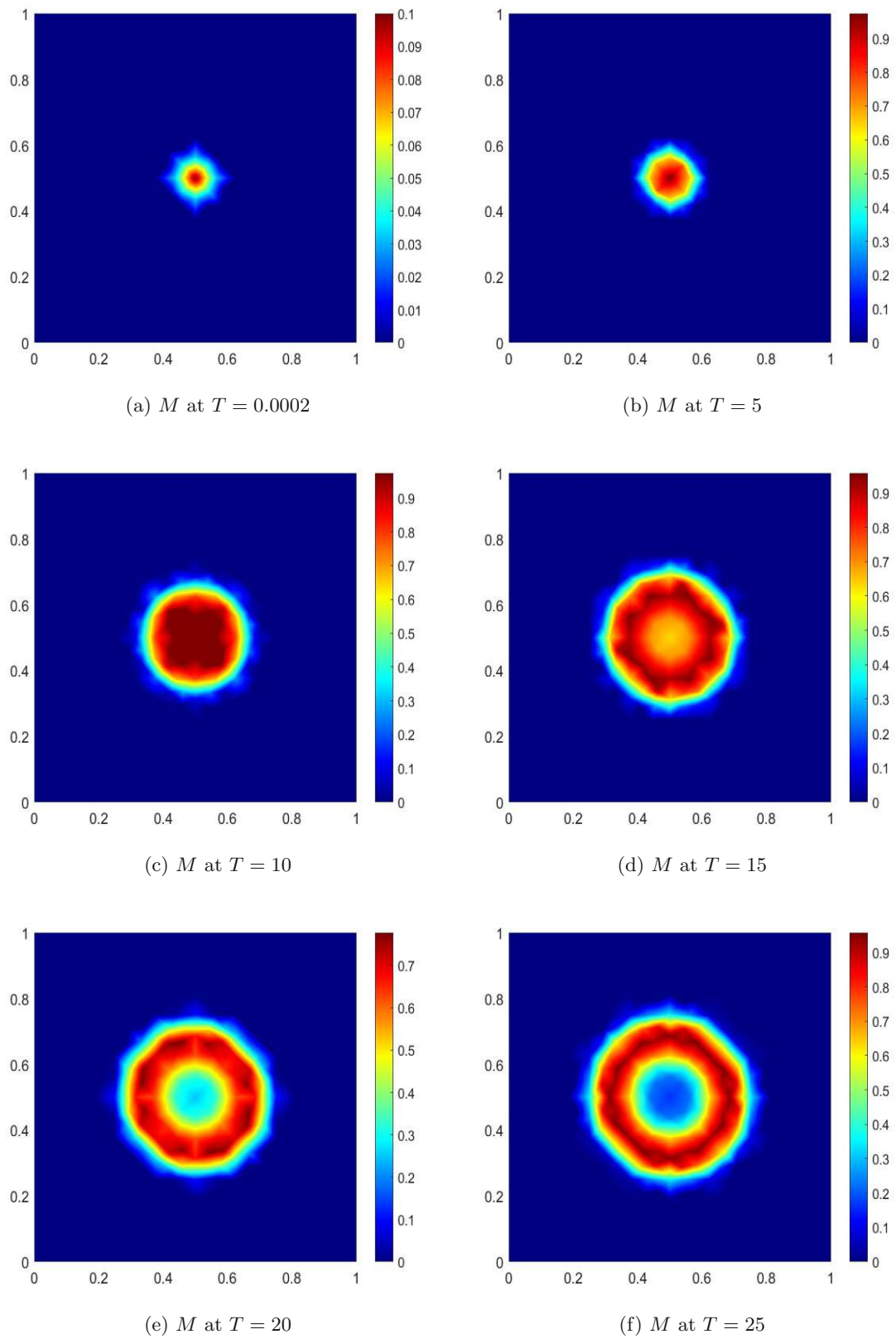
(a) $M$ at $T = 0.0002$

(b) $M$ at $T = 5$

(c) $M$ at $T = 10$

(d) $M$ at $T = 15$

(e) $M$ at $T = 20$

(f) $M$ at $T = 25$

Figure 3.6: Test case 3: Evolution of $M$.

(a) $N$ at $T = 0.0002$

(b) $N$ at $T = 5$

(c) $N$ at $T = 10$

(d) $N$ at $T = 15$

(e) $N$ at $T = 20$

(f) $N$ at $T = 25$

Figure 3.7: Test case 3: Evolution of $N$.

(a) $S$ at $T = 0.0002$

(b) $S$ at $T = 5$

(c) $S$ at $T = 10$

(d) $S$ at $T = 15$

(e) $S$ at $T = 20$

(f) $S$ at $T = 25$

Figure 3.8: Test case 3: Evolution of $S$.

(a) $A$ at $T = 0.0002$

(b) $A$ at $T = 5$

(c) $A$ at $T = 10$

(d) $A$ at $T = 15$

(e) $A$ at $T = 20$

(f) $A$ at $T = 25$

Figure 3.9: Test case 3: Evolution of $A$.

# 4 Existence Analysis for a Cahn-Hilliard-Type System Modeling Biofilm Growth

In this chapter, we provide the mathematical details for the model (1.24)–(1.28). To this end, we first introduce the main results and briefly describe the key ideas in section 4.1. Then, we introduce a truncated and regularized version of equations (1.24)–(1.28) in section 4.2.1 and prove the existence by using a galerkin approximation in section 4.2.2. In section 4.3, we provide uniform estimates in form of an energy and an entropy inequality which are necessary in order to perform the deregularization limit. In section 4.4, we conclude the existence proof by performing the limit. Finally, we present numerical experiments in section 4.5 where we also compare the numerics of the modified system (1.24)–(1.28) to the non–modified (yet simplified) system of [WZ12].

## 4.1 Main Result and Key Ideas

We impose the following assumptions:

(A1) Domain: $\Omega \subset \mathbb{R}^d$ $(d \geq 1)$ is a bounded domain with Lipschitz continuous boundary. Set $\Omega_T = \Omega \times (0, T)$.

(A2) Initial data: $u^0 \in H^1(\Omega)$ satisfies $0 < u_* \leq u^0 \leq u^*$ in $\Omega$ for some $u_*, u^* > 0$ and $v^0 \in L^2(\Omega)$ satisfying $0 \leq v^0 \leq 1$ in $\Omega$.

(A3) Source terms: $g_0 \in C^0([0,1])$ is nondecreasing and satisfies $g_0(0) = 0$, and $h_0 \in C^1([0,1])$ is nondecreasing.

Our main result is the global existence of bounded weak solutions.

**Theorem 22** (Global existence). *Let Assumptions (A1)–(A3) hold. Then there exists a weak solution $(u, v)$ to (1.24)–(1.28) with the constitutive relations (1.29)–(1.32), satisfying $0 \leq u \leq 1$, $0 \leq v \leq 1$ in $\Omega_T$,*

$$u \in L^2(0, T; H^2(\Omega)) \cap C^0([0, T]; H^1(\Omega)),$$
$$(1 - u)\nabla v, \ \partial_t u, \ \partial_t v \in L^2(0, T; H^1(\Omega)'),$$

*and the weak formulation for all $\phi_1, \phi_2 \in L^2(0, T; H^2(\Omega))$,*

$$\int_0^T \langle \partial_t v, \phi_1 \rangle dt + \int_0^T \langle (1 - u)\nabla v, \nabla \phi_1 \rangle dt = \int_0^T \int_\Omega g(u, v) \phi_1 \, dx \, dt,$$

$$\int_0^T \langle \partial_t u, \phi_2 \rangle dt + \int_0^T \langle J, \nabla \phi_2 \rangle dt = \int_0^T \int_\Omega h(u,v)\phi_2 dx dt,$$

*where $\langle \cdot, \cdot \rangle$ is the dual product between $H^1(\Omega)'$ and $H^1(\Omega)$ and $J = -\nabla(M(u)\Delta u) + M'(u)\nabla u + M(u)f''(u)\nabla u$ as well as $(1-u)\nabla v = \nabla((1-u)v) + v\nabla u$ are understood in the sense of $L^2(0,T;H^1(\Omega)')$.*

The proof of Theorem 22 is based on a suitable approximation scheme, truncating the nonlinearities and using a Galerkin method similarly as in [EG96]. Uniform estimates are obtained from the energy and entropy equalities, proved in Lemma 28 for the sequence of approximate solutions,

$$\frac{d}{dt}\int_\Omega \left(\frac{1}{2}|\nabla u|^2 + f(u)\right)dx + \int_\Omega M(u)|\nabla \mu|^2 dx = \int_\Omega h(u,v)\mu dx, \tag{4.1}$$

$$\frac{d}{dt}\int_\Omega \Phi(u)dx + \int_\Omega \left((\Delta u)^2 + f''(u)|\nabla u|^2\right)dx = \int_\Omega h(u,v)\Phi'(u)dx, \tag{4.2}$$

where $\Phi$ is defined by $\Phi''(u) = 1/M(u)$ and $\Phi(1/2) = \Phi'(1/2) = 0$. This function can be interpreted as the thermodynamic entropy of the system, since a computation shows that, with $M(u)$ given by (1.29),

$$\Phi(u) = u\log u + (1-u)\log(1-u) + \log 2 \geq 0 \quad \text{for } 0 < u < 1.$$

Since $f''(u)|\nabla u|^2 \geq -2\lambda|\nabla u|^2$ for $0 < u < 1$, the corresponding integral in (4.2) can be bounded by Gronwall's lemma and the energy bound (4.1).

The difficulty is to estimate the right-hand sides of (4.1)–(4.2). The term $h(u,v)$ contains the factor $u(1-u)$ which cancels the singularity from $\Phi'(u)$, such that $\int_\Omega h(u,v)\Phi'(u)dx$ is bounded. For the other integral, we include the definition of $\mu$ and integrate by parts:

$$\int_\Omega h(u,v)\mu dx = \int_\Omega \left((1-2u)|\nabla u|^2 + u(1-u)h_0'(v)\nabla v \cdot \nabla u + h(u,v)f'(u)\right)dx.$$

The last term is bounded since $h(u,v)$ cancels the singularity of the potential $f'(u)$. The first term can be treated by Gronwall's lemma since $u$ is bounded. For the second term, we use Young's inequality:

$$\int_\Omega u(1-u)h_0'(v)\nabla v \cdot \nabla u dx \leq C\int_\Omega |\nabla u|^2 dx + C\int_\Omega (1-u)|\nabla v|^2 dx,$$

where we use the property $0 \leq v \leq 1$. The last integral can be absorbed by the energy bound for $v$:

$$\frac{1}{2}\frac{d}{dt}\int_\Omega v^2 dx + \int_\Omega (1-u)|\nabla v|^2 dx = \int_\Omega g(u,v)v dx \leq C. \tag{4.3}$$

There is another difficulty: Because of the degeneracy in the equation for $v$, we do not obtain an estimate for $\nabla v$ (see (4.3)) and therefore we cannot expect strong convergence for (a subsequence of) the approximate solutions $(v_\delta)$ with $\delta > 0$ being an approximation parameter, but only weak* convergence in $L^\infty(\Omega_T)$. Surprisingly, the weak convergence of $(v_\delta)$ is enough to pass to the limit $\delta \to 0$ in $(1-u_\delta)\nabla v_\delta$, since this expression can be

written as $\nabla((1 - u_\delta)v_\delta) + v_\delta \nabla u_\delta$, which converges weakly in the sense of distributions, since $(\nabla u_\delta)$ converges strongly (up to a subsequence). However, the weak convergence is not sufficient to perform the limit in the reaction rates. The idea is to use the duality of $H^1(\Omega)'$ and $H^1(\Omega)$ as well as a Minty–Browder trick. Indeed, since $h_0$ is nondecreasing, we have for $y \in C_0^\infty(\Omega_T)$,

$$0 \le \int_0^T \int_\Omega u_\delta(1 - u_\delta)(v_\delta - y)(h_0(v_\delta) - h(y))dxdt$$

$$= \int_0^T \left\langle v_\delta - y, u_\delta(1 - u_\delta)(h_0(v_\delta) - h_0(y)) \right\rangle dt.$$

(Observe that we need to truncate the factor $u_\delta(1 - u_\delta)$, since we cannot expect that $0 \le u_\delta \le 1$; see Section 4.2.1.) By the Aubin–Lions lemma, $v_\delta \to v$ strongly in $L^2(0, T; H^1(\Omega)')$ and $u_\delta \to u$ strongly in $L^2(0, T; H^1(\Omega))$. Hence, a computation shows that the limit $\delta \to 0$ in the previous inequality leads to

$$0 \le \int_0^T \left\langle v - y, u(1 - u)(h_1 - h(y)) \right\rangle dt,$$

where $h_1$ is the weak $L^2(\Omega_T)$-limit of $(h_0(v_\delta))$. A Minty–Browder argument, made precise in Lemma 33, shows that $h_1 = h_0(v)$, implying that $h(u_\delta, v_\delta) \rightharpoonup h(u, v)$ weakly in $L^2(\Omega_T)$.

## 4.2 Existence for the approximate system

### 4.2.1 Truncated regularized system

We truncate the functions $M(u)$, $f(u)$, and the source terms. Let $\delta > 0$ and set $[u]_+ = \max\{0, u\}$ and $[u]_+^1 = \min\{1, \max\{0, u\}\}$ for $u \in \mathbb{R}$. We introduce for $u \in \mathbb{R}$

$$M_\delta(u) = \begin{cases} M(\delta) & \text{if } u \le \delta, \\ M(u) & \text{if } \delta < u < 1 - \delta, \\ M(1 - \delta) & \text{if } u \ge 1 - \delta. \end{cases}$$

Then $M_\delta(u) \ge M(u)$ for $u \in \mathbb{R}$. Furthermore, we set

$$D_+(u) = [1 - u]_+^1.$$

We approximate the singular part $f_1(u) = N^{-1}u \log u + (1 - u) \log(1 - u)$ of the free energy by setting

$f_{1,\delta}(u) = f_1(u) \quad \text{if } \delta < u < 1 - \delta,$

$f_{1,\delta}(u) = f_1(\delta) + f_1'(\delta)(u - \delta) + \frac{1}{2}f_1''(\delta)(u - \delta)^2 \quad \text{if } u \le \delta,$

$f_{1,\delta}(u) = f_1(1 - \delta) + f_1'(1 - \delta)(u - (1 - \delta)) + \frac{1}{2}f_1''(1 - \delta)(u - (1 - \delta))^2 \quad \text{if } u \ge 1 - \delta.$

This means that

$$f_{1,\delta}''(u) = \begin{cases} f_1''(\delta) & \text{if } u \le \delta, \\ f_1''(u) & \text{if } \delta < u < 1 - \delta, \\ f_1''(1 - \delta) & \text{if } u \ge 1 - \delta. \end{cases} \tag{4.4}$$

The regular part $f_2(u) = \lambda u(1 - u)$ $(0 \leq u \leq 1)$ of the free energy is extended to $\mathbb{R}$ such that $|f_2(u)| \leq C$ for $u \in \mathbb{R}$. Furthermore, we set $f_\delta = f_{1,\delta} + f_2$, and this function is defined for all $u \in \mathbb{R}$. We also need to truncate the source terms:

$$g_+(u, v) = -[u]_+^1 g_0([v]_+^1), \quad h_+(u, v) = [u]_+[1 - u]_+ h_0([v]_+^1).$$

Finally, let $\kappa > 0$. We wish to find a solution to the truncated and regularized system

$$\partial_t v - \operatorname{div}(D_+(u)\nabla v) - \kappa \Delta v = g_+(u, v), \tag{4.5}$$

$$\partial_t u - \operatorname{div}(M_\delta(u)\nabla \mu) = h_+(u, v), \tag{4.6}$$

$$\mu = -\Delta u + f'_\delta(u) \quad \text{in } \Omega, \ t > 0, \tag{4.7}$$

subject to the initial conditions (1.28) and the Neumann boundary conditions

$$\nabla v \cdot \nu = \nabla \mu \cdot \nu = \nabla u \cdot \nu = 0 \quad \text{on } \partial\Omega, \ t > 0. \tag{4.8}$$

### 4.2.2 Galerkin approximation

To solve (4.5)–(4.8) with initial condition (1.28), we use the Galerkin method (as in [EG96]). Let $(\phi_\ell)_{\ell \in \mathbb{N}}$ be the orthonormal eigenfunctions of the Laplace operator with homogeneous Neumann boundary conditions. We can assume that $\lambda_1 = 0$ and $\phi_1 = \text{const}$. Let $L \in \mathbb{N}$. We wish to find solutions

$$v_L(x, t) = \sum_{\ell=1}^L A_\ell(t)\phi_\ell(x), \quad u_L(x, t) = \sum_{\ell=1}^L B_\ell(t)\phi_\ell(x), \quad \mu_L(x, t) = \sum_{\ell=1}^L C_\ell(t)\phi_\ell(x)$$

to the finite-dimensional system

$$\int_\Omega \partial_t v_L \phi \, dx = -\int_\Omega (D_+(u_L) + \kappa)\nabla v_L \cdot \nabla \phi \, dx + \int_\Omega g_+(u_L, v_L)\phi \, dx, \tag{4.9}$$

$$\int_\Omega \partial_t u_L \phi \, dx = -\int_\Omega M_\delta(u_L)\nabla \mu_L \cdot \nabla \phi \, dx + \int_\Omega h_+(u_L, v_L)\phi \, dx, \tag{4.10}$$

$$\int_\Omega \mu_L \phi \, dx = \int_\Omega \nabla u_L \cdot \nabla \phi \, dx + \int_\Omega f'_\delta(u_L)\phi \, dx \tag{4.11}$$

for all $\phi \in \text{span}(\phi_1, \ldots, \phi_L)$, with the initial conditions

$$v_L(0) = \sum_{\ell=1}^L (v^0, \phi_\ell)_{L^2(\Omega)}\phi_\ell \, dx, \quad u_L(0) = \sum_{\ell=1}^L (u^0, \phi_\ell)_{L^2(\Omega)}\phi_\ell \, dx.$$

This gives an initial-value problem for a system of ordinary differential equations for $(A_1, \ldots, A_L)$ and $(B_1, \ldots, B_L)$:

$$\partial_t A_\ell = -\int_\Omega ([u_L]_+^1 + \kappa)\nabla v_L \cdot \nabla \phi_\ell \, dx + \int_\Omega g_+(u_L, v_L)\phi_\ell \, dx,$$

$$\partial_t B_\ell = -\int_\Omega M_\delta(u_L)\nabla \mu_L \cdot \nabla \phi_\ell \, dx + \int_\Omega h_+(u_L, v_L)\phi_\ell \, dx,$$

$$C_\ell = \int_\Omega \nabla u_L \cdot \nabla \phi_\ell dx + \int_\Omega f'_\delta(u_L)\phi_\ell dx \quad \text{for } \ell = 1, \dots, L,$$

with the initial conditions $A_\ell(0) = (v^0, \phi_\ell)_{L^2(\Omega)}$ and $B_\ell(0) = (u^0, \phi_\ell)_{L^2(\Omega)}$. As the right-hand side of this system is continuous in $(A_1, \dots, A_L)$ and $(B_1, \dots, B_L)$, the Peano theorem ensures the existence of a local solution. To extend this solution globally, we prove some a priori estimates.

**Lemma 23** (Energy estimate for the Galerkin approximation). *There exists a constant $C(\delta) > 0$ independent of $L$ such that for all $t \in (0, T)$,*

$$\frac{1}{2}\|\nabla u_L(t)\|^2_{L^2(\Omega)} + \int_\Omega f_\delta(u_L(t))dx + \frac{1}{2}M(\delta)\int_0^t \|\nabla\mu_L\|^2_{L^2(\Omega)}ds$$
$$\leq \frac{1}{2}\|\nabla u_L(0)\|^2_{L^2(\Omega)} + \int_\Omega f_\delta(u_L(0))dx + C(\delta),$$

*and because of Assumption (A2), the right-hand side can be bounded independently of $L$.*

*Proof.* We choose $\phi = \mu_L$ in (4.10) and $\phi = \partial_t u_L$ in (4.11):

$$\int_\Omega \partial_t u_L \mu_L dx = -\int_\Omega M_\delta(u_L)|\nabla\mu_L|^2 dx + \int_\Omega h_+(u_L, v_L)\mu_L dx$$
$$\leq -M(\delta)\int_\Omega |\nabla\mu_L|^2 dx + C\|\mu_L\|_{L^1(\Omega)},$$

$$\int_\Omega \mu_L \partial_t u_L dx = \int_\Omega \nabla u_L \cdot \nabla\partial_t u_L dx + \int_\Omega f'_\delta(u_L)\partial_t u_L dx$$
$$= \frac{d}{dt}\left(\frac{1}{2}\int_\Omega |\nabla u_L|^2 dx + \int_\Omega f_\delta(u_L)dx\right),$$

since $|h_+(u_L, v_L)| \leq C$ because of our truncations. Here and in the following, $C > 0$ denotes a generic constant with values changing from line to line. Equating both expressions and integrating over $(0, t)$ gives

$$\frac{1}{2}\int_\Omega \|\nabla u_L(t)\|^2_{L^2(\Omega)} + \int_\Omega f_\delta(u_L(t))dx \leq \frac{1}{2}\int_\Omega |\nabla u_L(0)|^2 dx + \int_\Omega f_\delta(u_L(0))dx \quad (4.12)$$
$$- M(\delta)\int_0^t \int_\Omega |\nabla\mu_L|^2 dx ds + C\int_0^t \|\mu_L\|_{L^1(\Omega)}ds.$$

The choice $\phi_1 = \text{const.}$ in (4.11) shows that

$$\left|\int_\Omega \mu_L dx\right| \leq \int_\Omega |f'_\delta(u_L)|dx \leq C(\delta).$$

Set $\bar{\mu}_L = |\Omega|^{-1}\int_\Omega \mu_L dx$. By the Poincaré–Wirtinger inequality, the previous estimate provides a bound for the $L^2(\Omega)$ norm of $\mu_L$:

$$\|\mu_L\|_{L^2(\Omega)} \leq \|\mu_L - \bar{\mu}_L\|_{L^2(\Omega)} + \|\bar{\mu}_L\|_{L^2(\Omega)} \leq C_P\|\nabla\mu_L\|_{L^2(\Omega)} + C(\delta).$$

Applying Young's inequality, we have

$$\int_0^t \|\mu_L\|_{L^1(\Omega)}ds \le C(\Omega)\int_0^t \|\mu_L\|_{L^2(\Omega)}ds \le \frac{1}{2}M(\delta)\int_0^t \|\nabla\mu_L\|_{L^2(\Omega)}^2 ds + C(\delta,\Omega,T).$$

Inserting this estimate into (4.12) finishes the proof. $\qquad\square$

**Lemma 24** (Estimates for $u_L$ and $\mu_L$)**.** *There exists $C(\delta) > 0$ independent of $L$ such that*

$$\|u_L\|_{L^\infty(0,T;H^1(\Omega))} + \|\mu_L\|_{L^2(0,T;H^1(\Omega))} \le C(\delta).$$

*Proof.* The proof of Lemma 23 shows that $(\nabla\mu_L)$ and $(\mu_L)$ are bounded in $L^2(\Omega_T)$ and that $(\nabla u_L)$ is bounded in $L^\infty(0,T;L^2(\Omega))$. We choose $\phi_1 = \text{const.}$ in (4.10):

$$\frac{d}{dt}\int_\Omega u_L dx = \int_\Omega h_+(u_L,v_L)dx \le C(\Omega).$$

Consequently, $\int_\Omega u_L(t)dx$ is uniformly bounded, at least on finite time intervals. This allows us to apply the Poincaré–Wirtinger inequality to deduce an $L^2(\Omega)$ bound for $u_L(t)$ uniformly in time. $\qquad\square$

We also need a priori estimates for the substrate concentration.

**Lemma 25** (Estimates for $v_L$)**.** *There exists $C(v^0) > 0$ only depending on the initial datum $v^0$ such that*

$$\|v_L\|_{L^\infty(0,T;L^2(\Omega))} + \|D_+(u_L)^{1/2}\nabla v_L\|_{L^2(\Omega_T)} + \kappa^{1/2}\|\nabla v^L\|_{L^2(\Omega_T)} \le C(v^0).$$

*Proof.* We choose the test function $\phi = v_L$ in (4.9) and take into account that $g_0(0) = 0$:

$$\frac{1}{2}\frac{d}{dt}\int_\Omega v_L^2 dx + \int_\Omega (D_+(u_L)+\kappa)|\nabla v_L|^2 dx = -\int_\Omega [u_L]_+^1 g_0([v_L]_+^1)v_L dx \le 0.$$

An integration over $(0,T)$ yields the result. $\qquad\square$

The uniform estimates for $u_L$ in $L^\infty(0,T;H^1(\Omega))$ and $v_L$ in $L^\infty(0,T;L^2(\Omega))$ show that the coefficients $(A_\ell)$ and $(B_\ell)$ are bounded in $(0,T)$. Thus, we infer the global existence of solutions to the Galerkin system (4.9)–(4.11). To pass to the limit $L \to \infty$, we need an estimate for the time derivatives.

**Lemma 26** (Estimates for the time derivatives)**.** *There exist $C_1(\delta) > 0$ depending on $\delta$ and $C_2 > 0$ independent of $\delta$ such that*

$$\|\partial_t u_L\|_{L^2(0,T;H^1(\Omega)')} \le C_1(\delta), \quad \|\partial_t v_L\|_{L^2(0,T;H^1(\Omega)')} \le C_2.$$

*Proof.* Let $\psi \in L^2(0,T;H^1(\Omega))$ and let $\Pi_L\psi$ be the projection of $\psi$ on $\text{span}(\phi_1,\dots,\phi_L)$. We infer from (4.9) and the bounds of Lemma 25 that

$$\left|\int_0^T\int_\Omega \partial_t v_L \psi\,dxdt\right| = \left|\int_0^T\int_\Omega \partial_t v_L \Pi_L\psi\,dxdt\right|$$

$$\leq \int_0^T \|D_+(u_L)^{1/2}\|_{L^\infty(\Omega)}\|D_+(u_L)^{1/2}\nabla v_L\|_{L^2(\Omega)}\|\nabla \Pi_L \psi\|_{L^2(\Omega)}dt$$

$$+ \int_0^t \|g_+(u_L,v_L)\|_{L^2(\Omega)}\|\Pi_L\psi\|_{L^2(\Omega)}dt \leq C_2\|\psi\|_{L^2(0,T;H^1(\Omega))}.$$

Furthermore, using $M_\delta(u_L) \leq C_M$ and the bounds of Lemma 24,

$$\left| \int_0^T \int_\Omega \partial_t u_L \psi dx dt \right| = \left| \int_0^T \int_\Omega \partial_t u_L \Pi_L \psi dx dt \right|$$

$$\leq C_M \int_0^T \|\nabla \mu_L\|_{L^2(\Omega)}\|\nabla \Pi_L \psi\|_{L^2(\Omega)}dt$$

$$+ \int_0^T \|h_+(u_L,v_L)\|_{L^2(\Omega)}\|\Pi_L\psi\|_{L^2(\Omega)}dt \leq C_1(\delta)\|\psi\|_{L^2(0,T;H^1(\Omega))}.$$

This concludes the proof. □

The estimates of Lemmas 24–26 allow us to apply the Aubin–Lions lemma [Sim87, Corollary 4] to find subsequences (not relabeled) such that, as $L \to \infty$,

$$u_L \to u, \quad v_L \to v \quad \text{strongly in } L^2(\Omega_T),$$

$$u_L \rightharpoonup u, \quad v_L \rightharpoonup v, \quad \mu_L \rightharpoonup \mu \quad \text{weakly in } L^2(0,T;H^1(\Omega)),$$

$$\partial_t u_L \rightharpoonup \partial_t u, \quad \partial_t v_L \rightharpoonup \partial_t v \quad \text{weakly in } L^2(0,T;H^1(\Omega)').$$

Since $M_\delta$, $D_+$, $f'_\delta$, $g_+$, and $h_+$ are bounded functions, we have

$$M_\delta(u_L) \to M_\delta(u), \quad D_+(u_L) \to D_+(u), \quad f'_\delta(u_L) \to f'_\delta(u),$$

$$h_+(u_L,v_L) \to h_+(u,v), \quad g_+(u_L,v_L) \to g_+(u,v) \quad \text{strongly in } L^2(\Omega_T).$$

Thus, we can perform the limit $L \to \infty$ in the Galerkin system (4.9)–(4.11), which yields the existence of a solution $(u,v,\mu)$ to

$$\int_0^t \langle \partial_t u, \phi_1 \rangle ds = -\int_0^t \int_\Omega M_\delta(u)\nabla\mu \cdot \nabla\phi_1 dx ds + \int_0^t \int_\Omega h_+(u,v)\phi_1 dx ds, \tag{4.13}$$

$$\int_0^t \langle \partial_t v, \phi_2 \rangle ds = -\int_0^t \int_\Omega (D_+(u)+\kappa)\nabla v \cdot \nabla\phi_2 dx ds + \int_0^t \int_\Omega g_+(u,v)\phi_2 dx ds, \tag{4.14}$$

$$\int_0^t \int_\Omega \mu\phi_3 dx ds = \int_0^t \int_\Omega \nabla u \cdot \nabla\phi_3 dx ds + \int_0^t \int_\Omega f'_\delta(u)\phi_3 dx ds \tag{4.15}$$

for all $\phi_i \in L^2(0,T;H^1(\Omega))$, $i = 1,2,3$, and all $0 < t < T$, recalling that $\langle \cdot,\cdot \rangle$ is the dual product between $H^1(\Omega)'$ and $H^1(\Omega)$.

## 4.3 Uniform estimates

We need some estimates uniform in $\delta$ and $\kappa$ as well as lower and upper bounds to remove the truncation.

**Lemma 27** (Uniform estimates for $v$)**.** *There exists $C(v^0) > 0$ only depending on $v^0$ such that*

$$\|v\|_{L^\infty(0,T;L^2(\Omega))} + \|D_+(u)^{1/2}\nabla v\|_{L^2(\Omega_T)} + \kappa^{1/2}\|\nabla v\|_{L^2(\Omega_T)} \leq C(v^0).$$

*Furthermore, it holds that $0 \leq v(t) \leq 1$ in $\Omega$ for $0 < t < T$.*

Because of the lower and upper bounds for $v$, we can remove the truncation in $g_+(u,v) = -[u]_+^1 g_0(v)$ and $h_+(u,v) = [u]_+[1-u]_+ h_0(v)$.

*Proof.* We start with the lower and upper bounds for $v$. We use the test function $[v]_- = \min\{0,v\}$ in (4.14) and use the assumption $v(0) \geq 0$ in $\Omega$:

$$\int_\Omega [v(t)]_-^2 dx + \int_0^t \int_\Omega (D_+(u) + \kappa)|\nabla[v]_-|^2 dx ds = \int_0^t \int_\Omega g_+(u,v)[v]_- dx ds = 0,$$

since $g_0(0) = 0$ implies that $g_+(u,v)[v]_- = -[u]_+^1 g_0(0)[v]_- = 0$. This implies that $v(t) \geq 0$ in $\Omega$, $t > 0$. The property $v(t) \leq 1$ is proved in a similar way using the test function $[v-1]_+$ and the fact that $g_+(u,v)[v-1]_+ = -[u]_+^1 g_0(1)[v-1]_+ \leq 0$. The remaining estimates can be shown as in Lemma 25. $\qquad\square$

Next, we show some uniform estimates for $u$. For this, we introduce the entropy density

$$\Phi_\delta(u) = \int_{1/2}^u \int_{1/2}^s \frac{dr\,ds}{M_\delta(r)} \geq 0. \tag{4.16}$$

**Lemma 28** (Energy and entropy estimates)**.** *There exists $C(T) > 0$ independent of $\delta$ and $\kappa$ such that for all $t > 0$ and all sufficiently small $\delta > 0$,*

$$\sup_{0<t<T} \int_\Omega \left(\frac{1}{2}|\nabla u(t)|^2 dx + f_\delta(u(t))\right)dx + \int_0^T \int_\Omega M_\delta(u)|\nabla\mu|^2 dx ds \leq C(T), \tag{4.17}$$

$$\sup_{0<t<T} \int_\Omega \Phi_\delta(u(t))dx + \int_0^T \int_\Omega (\Delta u)^2 dx ds \leq C(T). \tag{4.18}$$

Since $f_\delta$ is bounded from below (by construction), the energy inequality provides uniform bounds for $u$.

*Proof.* We first prove the energy inequality and then the entropy inequality.

*Step 1: Energy inequality.* We know from Section 4.2.2 that $u \in L^\infty(0,T;H^1(\Omega))$ and $\mu \in L^2(0,T;H^1(\Omega))$. Then we infer from the boundedness of $f_\delta'$ that $\Delta u = f_\delta'(u) - \mu \in L^2(0,T;L^2(\Omega))$. By elliptic regularity theory, $u \in L^2(0,T;H^2(\Omega))$. Moreover, $\nabla\Delta u = f_\delta''(u)\nabla u - \nabla\mu \in L^2(\Omega_T)$, which implies that $u \in L^2(0,T;H^3(\Omega))$ (this regularity is not uniform in $(\delta,\kappa)$). Consequently, $\Delta u \in L^2(0,T;H^1(\Omega))$ and

$$0 = \int_0^t \langle \partial_t u, \mu + \Delta u - f_\delta'(u)\rangle ds$$

$$= \int_0^t \langle \partial_t u, \mu\rangle ds - \frac{1}{2}\int_\Omega (|\nabla u(t)|^2 - |\nabla u(0)|^2)ds - \int_\Omega (f_\delta(u(t)) - f_\delta(u(0)))ds.$$

On the other hand, we use $\phi_2 = \mu \in L^2(0,T;H^1(\Omega))$ as a test function in (4.13):

$$\int_0^t \langle \partial_t u, \mu \rangle ds + \int_0^t \int_\Omega M_\delta(u)|\nabla\mu|^2 dx ds = \int_0^t \int_\Omega h_+(u,v)\mu dx.$$

This shows that, using the definition of $\mu$,

$$\frac{1}{2}\int_\Omega |\nabla u(t)|^2 dx + \int_\Omega f_\delta(u(t))dx + \int_0^t \int_\Omega M_\delta(u)|\nabla\mu|^2 dx ds = \frac{1}{2}\int_\Omega |\nabla u^0|^2 dx \qquad (4.19)$$

$$+ \int_\Omega f_\delta(u^0)dx + \int_0^t \int_\Omega \nabla h_+(u,v)\cdot\nabla u dx ds + \int_0^t \int_\Omega h_+(u,v)f_\delta'(u)dx ds.$$

It remains to estimate the last two integrals. For the last but one integral, we insert the definition of $h_+(u,v)$ and apply Young's inequality:

$$\int_0^t \int_\Omega \nabla h_+(u,v)\cdot\nabla u dx ds$$

$$= \int_0^t \int_\Omega 1_{\{0<u<1\}}\big((1-2u)h_0(v)|\nabla u|^2 + u(1-u)h_0'(v)\nabla v\cdot\nabla u\big)dx ds$$

$$\leq C\int_0^t \int_\Omega |\nabla u|^2 dx ds + C\int_0^t \int_\Omega 1_{\{0<u<1\}}(1-u)|\nabla v|^2 dx ds$$

$$\leq C\int_0^t \int_\Omega |\nabla u|^2 dx ds + C,$$

where the last step follows from Lemma 27, and $C > 0$ denotes here and in the following a constant independent of $\delta$ and $\kappa$.

For the last integral in (4.19), we observe that the function $s \mapsto -s(1-s)(N^{-1}\log s - \log(1-s) + N^{-1} - 1)$ is bounded in $[0,1]$. We insert the definition of $f_\delta'(u)$ and distinguish three cases. First, let $u \leq \delta$. Then

$$h_+(u,v)f_\delta'(u) = [u]_+[1-u]_+\big(f_\delta'(\delta) + f_1''(\delta)(u-\delta) + \lambda(1-2u)\big)$$

$$= [u]_+[1-u]_+\left(\frac{1}{N}\log\delta - \log(1-\delta) + \frac{1}{N} - 1\right)$$

$$+ [u]_+[1-u]_+(u-\delta)\left(\frac{1}{N\delta} + \frac{1}{1-\delta}\right) + \lambda[u]_+[1-u]_+(1-2u)$$

$$\leq \delta(1-\delta)\left(|\log(1-\delta)| + \frac{1}{N}\right) + \lambda\delta(1-\delta) \leq C,$$

using $[u]_+[1-u]_+ \leq \delta(1-\delta)$ and $u - \delta \leq 0$. Second, let $\delta < u < 1-\delta$. We have

$$h_+(u,v)f_\delta'(u) = u(1-u)\left(\frac{1}{N}\log u - \log(1-u) + \frac{1}{N} - 1\right) + \lambda u(1-u)(1-2u) \leq C,$$

since $z \mapsto z\log z$ is bounded in $[0,1]$. Finally, let $u \geq 1-\delta$ (and $\delta \leq 1/2$). We obtain

$$h_+(u,v)f_\delta'(u) = [u]_+[1-u]_+\left(\frac{1}{N}\log(1-\delta) - \log\delta + \frac{1}{N} - 1\right)$$

$$+ [u]_+[1-u]_+(u-\delta)\left(\frac{1}{N(1-\delta)} + \frac{1}{\delta}\right) + \lambda[u]_+[1-u]_+(1-2u)$$

$$\leq \delta(1-\delta)\left(|\log\delta| + \frac{1}{N}\right) + (1-\delta)\delta(1-2\delta)\left(\frac{1}{N(1-\delta)} + \frac{1}{\delta}\right) + \lambda \leq C.$$

This proves that, for $0 < t < T$,

$$\int_0^t \int_\Omega h_+(u,v)f_\delta'(u)dxds \leq C(\Omega, T).$$

Therefore, we infer from (4.19) that

$$\frac{1}{2}\int_\Omega |\nabla u(t)|^2 dx + \int_\Omega f_\delta(u(t))dx + \int_0^t \int_\Omega M_\delta(u)|\nabla\mu|^2 dxds$$

$$= \frac{1}{2}\int_\Omega |\nabla u^0|^2 dx + \int_\Omega f_\delta(u^0)dx + C\int_0^t \int_\Omega |\nabla u|^2 dxds + C.$$

Since $u^0$ is strictly positive and bounded away from one, there exists $\delta_0 > 0$ such that $f_\delta(u^0) = f(u^0)$ for $0 < \delta \leq \delta_0$. An application of Gronwall's lemma shows (4.17).

*Step 2: Entropy inequality.* Because of the truncation, we have $\nabla\Phi_\delta'(u) = \nabla u/M_\delta(u) \in L^2(\Omega_T)$, where $\Phi_\delta$ is defined in (4.16). Thus, we can use $\phi_1 = \Phi_\delta'(u)$ as a test function in (4.13):

$$\int_\Omega \Phi_\delta(u(t))dx - \int_\Omega \Phi_\delta(u(0))dx = \int_0^t \langle \partial_t u, \Phi_\delta'(u)\rangle ds \qquad (4.20)$$

$$= -\int_0^t \int_\Omega M_\delta(u)\nabla\mu \cdot \nabla\Phi_\delta'(u)dxds + \int_0^t \int_\Omega h_+(u,v)\Phi_\delta'(u)dxds$$

$$\leq -\int_0^t \int_\Omega \nabla(-\Delta u + f_\delta'(u))\cdot\nabla u dxds + \int_0^t \int_\Omega [u]_+[1-u]_+h_0(v)\Phi_\delta'(u)dxds.$$

The first integral on the right-hand side can be written as

$$-\int_0^t \int_\Omega \nabla(-\Delta u + f_\delta'(u))\cdot\nabla u dxds = -\int_0^t \int_\Omega (\Delta u)^2 dxds - \int_0^t \int_\Omega f_\delta''(u)|\nabla u|^2 dxds.$$

Because of $f_{1,\delta}''(u) \geq 0$ by (4.4) and $f_2''(u) \geq -C$, we obtain

$$-\int_0^t \int_\Omega \nabla(-\Delta u + f_\delta'(u))\cdot\nabla u dxds \leq -\int_0^t \int_\Omega (\Delta u)^2 dxds + C\int_0^t \int_\Omega |\nabla u|^2 dxds.$$

We claim that the integrand of the last integral in (4.20) is bounded, i.e. $[u]_+[1-u]_+\Phi_\delta'(u)$ is bounded uniformly in $u \in [0,1]$ and $\delta \in (0,1/2)$. Indeed, if $\delta \leq u \leq 1-\delta$, we can compute

$$|[u]_+[1-u]_+\Phi_\delta'(u)| = \left|u(1-u)\int_{1/2}^u \frac{ds}{s(1-s)}\right| = \left|u(1-u)\log\frac{u}{1-u}\right| \leq 1.$$

If $0 < u < \delta$, we find that

$$|[u]_+[1-u]_+\Phi'_\delta(u)| = \left| u(1-u)\left( \int_{1/2}^\delta \frac{ds}{s(1-s)} + \int_\delta^u \frac{ds}{\delta(1-\delta)} \right) \right|$$
$$= u(1-u)\log\frac{\delta}{1-\delta} + u(1-u)\frac{\delta-u}{\delta(1-\delta)}.$$

The first term is uniformly bounded since $|u\log\delta| \leq |\delta\log\delta| \leq 1$ and $|(1-u)\log(1-\delta)| \leq 1$. This holds also true for the second term because of $u(1-u) < \delta(1-\delta)$. The final case $1-\delta < u < 1$ is treated in a similar way:

$$|[u]_+[1-u]_+\Phi'_\delta(u)| = \left| u(1-u)\left( \int_{1/2}^{1-\delta} \frac{ds}{s(1-s)} + \int_{1-\delta}^u \frac{ds}{\delta(1-\delta)} \right) \right|$$
$$= u(1-u)\log\frac{1-\delta}{\delta} + u(1-u)\frac{u-(1-\delta)}{\delta(1-\delta)}.$$

The first term is uniformly bounded since $|(1-u)\log\delta| \leq |\delta\log\delta| \leq 1$ and $|u\log(1-\delta)| \leq 1$, and the second term is bounded too. We conclude from (4.20) that

$$\int_\Omega \Phi_\delta(u(t))dx + \int_0^t \int_\Omega (\Delta u)^2 dx\,ds \leq \int_\Omega \Phi_\delta(u^0)dx + C\int_0^t \int_\Omega |\nabla u|^2 dx\,ds,$$

and the energy bound (4.17) leads to (4.18). $\qquad\square$

Finally, we derive a bound for the time derivatives of $u$ and $v$.

**Lemma 29** (Bounds for the time derivatives)**.** *There exists $C > 0$ independent of $\delta$ and $\kappa$ such that*

$$\|\partial_t u\|_{L^2(0,T;H^1(\Omega)')} + \|\partial_t v\|_{L^2(0,T;H^1(\Omega)')} \leq C.$$

*Proof.* The proof is similar to that one of Lemma 26; we just have to estimate the reaction terms. Since $0 \leq v \leq 1$, we have the pointwise bounds $g_+(u,v) = -[u]_+^1 g_0(v) \leq \max_{0 \leq v \leq 1} g_0(v)$ and $h_+(u,v) = [u]_+[1-u]_+ h_0(v) \leq \max_{0 \leq v \leq 1} h_0(v)$. Consequently, $\|g_+(u, v)\|_{L^2(\Omega_T)}$ and $\|h_+(u,v)\|_{L^2(\Omega_T)}$ are uniformly bounded, concluding the proof. $\qquad\square$

## 4.4 The limit $(\delta, \kappa) \to 0$

Set $\kappa = \delta$ and let $(u_\delta, v_\delta, \mu_\delta)$ be a weak solution to (4.13)–(4.15). Lemmas 27–29 give the following uniform bounds:

$$0 \leq v_\delta \leq 1 \quad \text{in } \Omega_T,$$
$$\|D_+(u_\delta)^{1/2}\nabla v_\delta\|_{L^2(\Omega_T)} + \delta^{1/2}\|v_\delta\|_{L^2(0,T;H^1(\Omega))} + \|\partial_t v_\delta\|_{L^2(0,T;H^1(\Omega)')} \leq C,$$
$$\|u_\delta\|_{L^\infty(0,T;H^1(\Omega))} + \|u_\delta\|_{L^2(0,T;H^2(\Omega))} + \|\partial_t u_\delta\|_{L^2(0,T;H^1(\Omega)')} \leq C,$$
$$\|M_\delta(u_\delta)^{1/2}\nabla \mu_\delta\|_{L^2(\Omega_T)} \leq C.$$

The Aubin–Lions lemma [Sim87, Corollary 4] implies the existence of a subsequence, which is not relabeled, such that, as $\delta \to 0$,

$$u_\delta \to u \quad \text{strongly in } L^2(0, T; H^1(\Omega)) \text{ and } C^0([0, T]; L^2(\Omega)).$$

We also have the weak convergences

$$v_\delta \rightharpoonup v \quad \text{weakly* in } L^\infty(0, T; L^\infty(\Omega)),$$
$$\partial_t u_\delta \rightharpoonup \partial_t u, \quad \partial_t v_\delta \rightharpoonup \partial_t v \quad \text{weakly in } L^2(0, T; H^1(\Omega)'),$$
$$D_+(u_\delta)\nabla v_\delta \rightharpoonup I, \quad M_\delta(u_\delta)^{1/2}\nabla\mu_\delta \rightharpoonup J \quad \text{weakly in } L^2(\Omega_T),$$

where $I, J \in L^2(\Omega_T)$, and it holds that $\delta\nabla v_\delta \to 0$ strongly in $L^2(\Omega_T)$. Before we identify the limits $I$ and $J$, we show that the limit $u$ is bounded from below and above.

**Lemma 30** ($L^\infty$ bounds for $u$). *It holds that $0 \leq u \leq 1$ in $\Omega_T$.*

*Proof.* We proceed as in the proofs of [EG96, Lemma 2] or [PP21, Theorem 5]. Let $\alpha > 0$ and introduce the set $V_{\alpha,\delta} = \{(x, t) \in \Omega_T : u_\delta(x, t) \geq 1 + \alpha\}$. Integrating $\Phi_\delta''(u_\delta(x, t)) = 1/M_\delta(1 - \delta) = 1/(\delta(1 - \delta))$ for $(x, t) \in V_{\alpha,\delta}$ twice gives

$$\Phi_\delta(u_\delta(x, t)) = \int_{1/2}^{u_\delta(x,t)} \int_{1/2}^{s} \frac{dr\,ds}{M_\delta(r)} = \frac{(u_\delta - 1/2)^2}{2\delta(1 - \delta)} \quad \text{for } (x, t) \in V_{\alpha,\delta}.$$

The entropy estimate (4.18) shows that

$$\frac{\alpha^2 |V_{\alpha,\delta}|}{2\delta(1 - \delta)} \leq \int_{V_{\alpha,\delta}} \frac{(u_\delta - 1/2)^2}{2\delta(1 - \delta)} d(x, t) = \int_{V_{\alpha,\delta}} \Phi_\delta(u_\delta) d(x, t) \leq C(T).$$

Then we deduce from the a.e. pointwise limit $u_\delta(x, t) \to u(x, t)$ as $\delta \to 0$ amd Fatou's lemma that

$$|\{u(x, t) \geq 1 + \alpha\}| = \lim_{\delta\to 0} |V_{\alpha,\delta}| \leq \lim_{\delta\to 0} \frac{2C(T)}{\alpha^2}\delta(1 - \delta) = 0,$$

implying that $u(x, t) \leq 1 + \alpha$ a.e. in $\Omega_T$ for all $\alpha > 0$. Therefore, $u(x, t) \leq 1$ in $\Omega_T$.

A similar argument proves that $u \geq 0$ in $\Omega_T$. Indeed, let $W_{\alpha,\delta} = \{(x, t) : u_\delta(x, t) \leq -\alpha\}$ for $\alpha > 0$. It follows from $\Phi_\delta''(u_\delta(x, t)) = 1/\delta(1 - \delta)$ for $(x, t) \in W_{\alpha,\delta}$ that $\Phi_\delta(u_\delta(x, t)) \leq (1/2 - u_\delta(x, t))^2/(2\delta(1 - \delta))$. Hence,

$$\frac{\alpha^2 |W_{\alpha,\delta}|}{2\delta(1 - \delta)} \leq \int_{W_{\alpha,\delta}} \frac{(1/2 - u_\delta)^2}{2\delta(1 - \delta)} d(x, t) = \int_{W_{\alpha,\delta}} \Phi_\delta(u_\delta) d(x, t) \leq C(T),$$

and proceeding as before gives $|\{u(x, t) \leq -\alpha\}| = 0$ in the limit $\delta \to 0$ for all $\alpha > 0$ and therefore $u \geq 0$ in $\Omega_T$. $\qquad\square$

We continue by identifying $I$. We conclude from $[1 - u_\delta]_+^1 v_\delta \rightharpoonup (1 - u)v$ and $v_\delta\nabla u_\delta \rightharpoonup v\nabla u$ weakly in $L^2(\Omega_T)$ that

$$D_+(u_\delta)\nabla v_\delta = \nabla([1 - u_\delta]_+^1 v_\delta) + v_\delta 1_{\{0<u_\delta<1\}}\nabla u_\delta \rightharpoonup \nabla((1 - u)v) + v\nabla u = (1 - u)\nabla v$$

weakly in $L^2(0, T; H^1(\Omega)')$. This shows that $I = (1 - u)\nabla v$ in $L^2(0, T; H^1(\Omega)')$.

**Lemma 31** (Identification of $J$). *It holds that* $J = -\nabla(M(u)\Delta u) + \nabla M(u)\Delta u + M(u)\nabla f'(u)$ *in the sense of* $L^2(0, T; H^1(\Omega)')$.

*Proof.* We proceed as in [EG96, Section 3]. It holds for $\phi \in C_0^\infty(\Omega_T)$ that

$$\int_0^T \int_\Omega M_\delta(u_\delta)\nabla\mu_\delta \cdot \nabla\phi \, dx dt = \int_0^T \int_\Omega M_\delta(u_\delta)\nabla\big(-\Delta u_\delta + f_\delta'(u_\delta)\big) \cdot \nabla\phi \, dx dt$$

$$= \int_0^T \int_\Omega M_\delta(u_\delta)\Delta u_\delta\Delta\phi \, dx dt + \int_0^T \int_\Omega M_\delta'(u_\delta)\Delta u_\delta\nabla u_\delta \cdot \nabla\phi \, dx dt$$

$$+ \int_0^T \int_\Omega M_\delta(u_\delta)f_\delta''(u_\delta)\nabla u_\delta \cdot \nabla\phi \, dx dt =: J_1 + J_2 + J_3.$$

First, we consider $J_1$. We observe that $M_\delta \to M$ uniformly, since by the mean-value theorem,

$$|M_\delta(z) - M(z)| \le \sup_{0 < z < \delta} |M(\delta) - M(z)| + \sup_{1-\delta < z < 1} |M(1-\delta) - M(z)|$$

$$\le M'(\xi_\delta)\delta + M'(\eta_\delta)\delta \to 0,$$

where $\xi_\delta \in (z, \delta)$ and $\eta_\delta \in (1 - \delta, z)$. This implies that $M_\delta(u_\delta) \to M(u)$ a.e. in $\Omega_T$ and, as $M_\delta$ is uniformly bounded, also strongly in $L^2(\Omega_T)$. Together with the convergence $\Delta u_\delta \rightharpoonup \Delta u$ weakly in $L^2(\Omega_T)$, we find that

$$J_1 \to \int_0^T \int_\Omega M(u)\Delta u\Delta\phi \, dx dt.$$

For the integral $J_2$, we claim that $M_\delta'(u_\delta)\nabla u_\delta \to M'(u)\nabla u$ strongly in $L^2(\Omega_T)$. This limit is not trivial since $M_\delta'$ is discontinuous at $\delta$ and $1 - \delta$. We consider the integrals

$$\int_0^T \int_\Omega |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt = \int_0^T \int_{\{0 < u < 1\}} |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt$$

$$+ \int_0^T \int_{\{u=0\}} |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt + \int_0^T \int_{\{u=1\}} |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt.$$

On the set $\{0 < u < 1\}$, we know that $M_\delta'(u_\delta) \to M'(u)$ a.e. in $\Omega_T$ and, because of the strong convergence of $(\nabla u_\delta)$, also $M_\delta'(u_\delta)\nabla u_\delta \to M'(u)\nabla u$ a.e. in $\Omega_T$ (possibly for a subsequence). Moreover, $|M_\delta'(u_\delta)\nabla u_\delta|^2$ is uniformly bounded on $\{0 < u < 1\}$. Therefore, by dominated convergence,

$$\int_0^T \int_{\{0 < u < 1\}} |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt \to 0.$$

It follows from $\nabla u = 0$ on $\{u = 0\} \cup \{u = 1\}$ and the uniform bound for $M_\delta'$ that

$$\int_0^T \int_{\{u=0\}} |M_\delta'(u_\delta)\nabla u_\delta - M'(u)\nabla u|^2 dx dt = \int_0^T \int_{\{u=0\}} |M_\delta'(u_\delta)\nabla u_\delta|^2 dx dt$$

$$\leq C \int_0^T \int_{\{u=0\}} |\nabla u_\delta|^2 dxdt \to \int_0^T \int_{\{u=0\}} |\nabla u|^2 dxdt = 0.$$

The limit in the remaining integral over $\{u = 1\}$ vanishes in the same way. This shows that

$$J_2 \to \int_0^T \int_\Omega M'(u)\Delta u \nabla u \cdot \nabla \phi dxdt.$$

Finally, for the limit in $J_3$, we observe that $M_\delta(z)f_\delta''(z) = M_\delta(z)(f_{1,\delta}''(z) + f_2''(z))$ is uniformly bounded, since the singularities as $\delta \to 0$ in $f_{1,\delta}''$ are canceled by the factor $M_\delta(z)$. Thus, it remains to show that $M_\delta(u_\delta)f_\delta''(u_\delta) \to M(u)f''(u)$ in $\Omega_T \setminus N$, where $N$ is a set of measure zero. To this end, we distinguish several cases.

Let $(x, t) \in \Omega_T \setminus N$ and $0 < u(x, t) < 1$. For given $\varepsilon > 0$, there exists $0 < \delta < \varepsilon$ such that $\delta < \varepsilon \leq u_\delta(x,t) \leq 1 - \varepsilon < 1 - \delta$. At this point, we have $M_\delta(u_\delta(x,t))f_\delta''(u_\delta(x,t)) = M(u_\delta(x,t))f''(u_\delta(x,t)) \to M(u(x,t))f''(u(x,t))$. Next, if $u(x,t) = 1$, we choose $\delta > 0$ such that $u_\delta(x,t) \geq 1 - \delta$. Then

$$M_\delta(u_\delta(x,t))f_\delta''(u_\delta(x,t)) = M(\delta)(f_1''(\delta) + f_2(u_\delta))$$
$$= N^{-1}\delta + (1 - \delta) + \delta(1 - \delta)f_2(u_\delta) \to 1 = (Mf'')(1).$$

On the other hand, if $u_\delta(x,t) < 1 - \delta$ and $u_\delta(x,t) \to 1$,

$$M_\delta(u_\delta(x,t))f_\delta''(u_\delta(x,t)) = M(u_\delta(x,t))f''(u_\delta(x,t))$$
$$= N^{-1}(1 - u_\delta(x,t)) + u_\delta(x,t) + u_\delta(1 - u_\delta)f_(''u_\delta) \to 1 = (Mf'')(1).$$

The case $u(x,t) = 0$ is treated in a similar way. We conclude that $M_\delta(u_\delta)f_\delta''(u_\delta) \to M(u)f''(u)$ strongly in $L^2(\Omega_T)$. Then, in view of the strong convergence of $(\nabla u_\delta)$,

$$J_3 \to \int_0^T \int_\Omega M(u)f''(u)\nabla u \cdot \nabla \phi dxdt.$$

Summarizing, we have shown that

$$\int_0^T \int_\Omega M_\delta(u_\delta)\nabla \mu_\delta \cdot \nabla \phi dxdt \to \int_0^T \int_\Omega \big(M(u)\Delta u \Delta \phi + M'(u)\Delta u \nabla u \cdot \nabla \phi$$
$$+ M(u)f''(u)\nabla u \cdot \nabla \phi\big)dxdt,$$

and the right-hand side can be identified as the weak formulation of $J$. □

**Remark 32.** *Choosing the mobility such that $\Phi(0) = \Phi(1) = \infty$, one can show that $\{u = 0\} \cup \{u = 1\}$ has measure zero, which means that $0 < u < 1$ holds a.e. in $\Omega_T$, and we can write $J = M(u)\nabla(-\Delta u + f'(u))$ in the sense of distributions. The claim that $\{u = 0\} \cup \{u = 1\}$ has measure zero can be proved as in [EG96, Corollary]. It follows from the entropy bound $\int_\Omega \Phi_\delta(u_\delta(t))dx \leq C(T)$ and the fact that $\liminf_{\delta \to 0} \Phi_\delta(u_\delta) = \Phi(u)$ if $0 < u < 1$ and $\liminf_{\delta \to 0} \Phi_\delta(u_\delta) = \infty$ else.*

It remains to pass to the limit $\delta \to 0$ in the reaction terms. Since $(v_\delta)$ is only converging weakly, this limit is not trivial. The idea is to use the Minty–Browder trick, which is possible since $(u_\delta)$ converges strongly in $L^2(0, T; H^1(\Omega))$.

**Lemma 33.** *It holds that $g_+(u_\delta, v_\delta) \rightharpoonup g(u,v)$ and $h_+(u_\delta, v_\delta) \rightharpoonup h(u,v)$ weakly in $L^2(\Omega_T)$ as $\delta \to 0$.*

*Proof.* We only show the limit in $h_+(u_\delta, v_\delta)$ as the proof in $g_+(u_\delta, v_\delta)$ is similar. We know that $(\partial_t v_\delta)$ is bounded in $L^2(0,T;H^1(\Omega)')$ and $(v_\delta)$ is bounded in $L^2(\Omega_T)$. Since the embedding $L^2(\Omega) \hookrightarrow H^1(\Omega)'$ is compact, we infer from the Aubin–Lions lemma that, up to a subsequence, $v_\delta \to v$ strongly in $L^2(0,T;H^1(\Omega)')$. Moreover, $([1 - u_\delta]_+^{1/2} \nabla v_\delta)$ is bounded in $L^2(\Omega_T)$. Furthermore, we know that $(u_\delta)$ is bounded in $L^\infty(0,T;H^1(\Omega))$ and $L^2(0,T;H^2(\Omega))$, and $u_\delta \to u$ strongly in $L^2(0,T;H^1(\Omega))$.

Let $y \in C_0^\infty(\Omega_T)$. It follows from the monotonicity of $h_0$ that

$$0 \le \int_0^T \int_\Omega [u_\delta]_+ [1 - u_\delta]_+ (v_\delta - y)(h_0(v_\delta) - h_0(y))dxdt \qquad (4.21)$$
$$= \int_0^T \big\langle v_\delta - y, [u_\delta]_+ [1 - u_\delta]_+ (h_0(v_\delta) - h_0(y)) \big\rangle dt,$$

recalling that $\langle \cdot, \cdot \rangle$ is the dual product between $H^1(\Omega)'$ and $H^1(\Omega)$. This formulation is possible if $[u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta) \in L^2(0,T;H^1(\Omega))$. To verify this statement, we observe that $\nabla u_\delta \in L^2(0,T;H^1(\Omega))$ implies that $(1 - 2u_\delta)\mathbf{1}_{\{0 < u_\delta < 1\}} \nabla u_\delta \in L^2(0,T;L^2(\Omega))$. Moreover, $[u_\delta]_+ [1 - u_\delta]_+^{1/2} \in L^\infty(\Omega_T)$ and $[1 - u_\delta]_+^{1/2} \nabla v_\delta \in L^2(\Omega_T)$. This shows that

$$\nabla\big([u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta)\big) = [u_\delta]_+ [1 - u_\delta]_+ h_0'(v_\delta) \nabla v_\delta + h_0(v_\delta)(1 - 2u_\delta)\mathbf{1}_{\{0 < u_\delta < 1\}} \nabla u_\delta$$

is a function in $L^2(\Omega_T)$, so that $[u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta) \in L^2(0,T;H^1(\Omega))$.

Let $h_1$ be the weak* limit of $(h_0(v_\delta))$ in $L^\infty(0,T;L^\infty(\Omega))$ and $h_2$ be the weak limit of $([u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta))$ in $L^2(\Omega_T)$. We claim that $h_2 = u(1 - u)h_1$. Indeed, since $(u_\delta)$ converges strongly in $L^2(0,T;H^1(\Omega))$, $[u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta) \rightharpoonup u(1-u)h_1$ weakly in $L^2(\Omega_T)$ (here, we use $0 \le u \le 1$ in $\Omega_T$; see Lemma 30), and we deduce from the uniqueness of the limit that $u(1 - u)h_1 = h_2$.

We can now pass to the limit $\delta \to 0$ in (4.21) to find that

$$0 \le \int_0^T \big\langle v - y, u(1-u)(h_1 - h_0(y)) \big\rangle dt = \int_0^T \int_\Omega u(1-u)(h_1 - h_0(y))(v - y)dxdt.$$

By density, this inequality holds for all $y \in L^2(\Omega_T)$. Let $w \in L^2(\Omega_T)$ and choose $y = v - \eta w$ for $\eta \in \mathbb{R}$. Then

$$0 \le \eta \int_0^T \int_\Omega u(1-u)(h_1 - h_0(v - \eta w))w\,dxdt.$$

Choosing $\eta > 0$ and performing the limit $\eta \to 0$ yields $\int_0^T \int_\Omega u(1-u)(h_1 - h_0(v))w\,dxdt \ge 0$. On the other hand, if $\eta < 0$ and $\eta \to 0$, we have $\int_0^T \int_\Omega u(1-u)(h_1 - h_0(v))w\,dxdt \le 0$. Since $w$ is arbitrary, $u(1-u)h_1 = u(1-u)h_0(v)$. Thus,

$$h_+(u_\delta, v_\delta) = [u_\delta]_+ [1 - u_\delta]_+ h_0(v_\delta) \rightharpoonup u(1-u)h_0(v) \quad \text{weakly in } L^2(\Omega_T).$$

This ends the proof. $\qquad\qquad\qquad\square$

**Remark 34** (Generalizations). *It is possible to generalize the relations* (1.29) *and* (1.32) *for the mobility and the reaction rates. For instance, we may choose* $M(u) = u^m(1-u)^m M_0(u)$ *for* $m \geq 1$ *and* $0 < m_* \leq M(u) \leq m^*$ *for* $u \in [0,1]$, *where* $m^* \geq m_* > 0$; *see [EG96]. In fact, we just need* $M(0) = M(1) = 0$ *and* $M(u)f''(u) \in C^0([0,1])$; *see [PP21]. The latter condition is needed to identify the weak limit J. The reaction terms may be generalized to* $g(u,v) = g_0(v)g_1(u)$ *and* $h(u,v) = h_0(v)h_1(u)$, *for instance, where we assume that* $g_1$ *is bounded in* $[0,1]$; $g_0$ *grows at most linearly;* $h_1$ *satisfies* $h_1(u)f'(u) \leq C$ *for all* $u[0,1]$ *to cancel the singularities of* $f'$; *and* $|h_1(u)| \leq C(1-u)$ *for* $u \in [0,1]$ *to estimate in Step 1 of the proof of Lemma 28 the integral*

$$\int_\Omega h_1(u)h_0'(v)\nabla v \cdot \nabla u \, dx \leq \int_\Omega |\nabla u|^2 dx + C \int_\Omega (1-u)|\nabla v|^2 dx.$$

*Clearly, also the free energy* $f(u)$ *may be generalized if the factors in the diffusion and reaction terms are adapted in such a way that the singularities from* $f'(u)$ *are canceled.*

## 4.5 Numerical experiments

### 4.5.1 Scaling of the equations

The biofilm model with physical units reads as follows:

$$\partial_t v - \text{div}(D(1-u)\nabla v) = -R_c uv,$$

$$\partial_t u - \text{div}(M'u(1-u)\nabla\mu) = u(1-u)\frac{R_p v}{K_v + v},$$

$$\mu = -\Gamma_1 \Delta u + \Gamma_2 f'(u),$$

and $f'(u)$ is given by (1.31), observing that the parameters $N$ and $\lambda$ and the volume fraction $u$ are dimensionless. Here, $D > 0$ is the diffusivity, $M' > 0$ the mobility constant, $R_c > 0$ the consumption rate, $R_p > 0$ the production rate, $\Gamma_1 > 0$ the parameter of the distortional energy, and $\Gamma_2 > 0$ the parameter of the mixing free energy.

Choosing the characteristic length $x_0$, the characteristic time $t_0$, the characteristic concentration $v_0$, and the characteristic chemical potential $\mu_0$, the scaled equations read as follows:

$$\partial_t v - \text{div}(D_0(1-u)\nabla v) = -R_c^0 uv, \tag{4.22}$$

$$\partial_t u - \text{div}(M_0 u(1-u)\nabla\mu) = u(1-u)\frac{R_p^0 v}{K + v}, \tag{4.23}$$

$$\mu = -\Gamma_1^0 \Delta u + \Gamma_2^0 f'(u), \tag{4.24}$$

where the dimensionless parameters are

$$D_0 = \frac{Dt_0}{x_0^2}, \quad M_0 = \frac{M't_0\mu_0}{x_0^2}, \quad R_c^0 = R_c t_0, \quad R_p^0 = R_p t_0,$$

$$K = \frac{K_v}{v_0}, \quad \Gamma_1^0 = \frac{\Gamma_1}{\mu_0 x_0^2}, \quad \Gamma_2^0 = \frac{\Gamma_2}{\mu_0}.$$

The model of [WZ12] (without elastic energy contributions) reads as

$$\partial_t((1-u)v) - \text{div}(D_0(1-u)\nabla v) = -u\frac{R_c v}{\widetilde{K}+v},$$

$$\partial_t u - \text{div}(M_0(1-u)\nabla\mu) = u\frac{R_p v}{K_v + v},$$

$$\mu = -\Gamma_1\Delta u + \Gamma_2 f'(u).$$

| Symbol | Parameter | Value | Unit |
|---|---|---|---|
| $D$ | Diffusivity | $10^{-10}$ | $\text{m}^2\,\text{s}^{-1}$ |
| $M'$ | Mobility | $2.5\cdot 10^{-8}$ | s |
| $R_c$ | Consumption rate | $10^{-2}$ | $\text{s}^{-1}$ |
| $R_p$ | Production rate | $10^{-2}$ | $\text{kg}\,\text{m}^{-3}\,\text{s}^{-1}$ |
| $K_v$ | Half-saturation constant | $10^{-4}$ | $\text{kg}\,\text{m}^{-3}$ |
| $\Gamma_1$ | Distortional energy | $4\cdot 10^{-15}$ | $\text{m}^4\,\text{s}^{-2}$ |
| $\Gamma_2$ | Mixing free energy | $4\cdot 10^{-6}$ | $\text{m}^2\,\text{s}^{-2}$ |
| $N$ | Polymerization parameter | $10^3$ | |
| $\lambda$ | Flory–Huggins parameter | $0.55$ | |
| $x_0$ | Characteristic length | $10^{-4}$ | m |
| $t_0$ | Characteristic time | $10^2$ | s |
| $v_0$ | Characteristic concentration | $10^{-3}$ | $\text{kg}\,\text{m}^{-3}$ |
| $k_B T$ | Thermal energy at $T = 300\,\text{K}$ | $4\cdot 10^{-21}$ | $\text{kg}\,\text{m}^2\,\text{s}^{-2}$ |
| $\widetilde{K}$ | Half-saturation constant for model of [WZ12] | $5\cdot 10^{-4}$ | |

Table 4.1: Parameters used in the numerical simulations.

The characteristic chemical potential $\mu_0$ is determined by the thermal energy and the characteristic concentration and length (see Table 4.1) as $\mu_0 = k_B T/(v_0 x_0^3) = 4\cdot 10^{-6}\,\text{m}^2\text{s}^{-2}$. The values of the physical parameters in Table 4.1 differ from those in [WZ12] but are of a similar order. With our values, the scaled parameters are of order one (except $K$ and $\Gamma_1^0$):

$$D_0 = R_c^0 = R_p^0 = 1, \quad K = 10^{-1}, \quad M_0 = 10^{-3}, \quad \Gamma_1^0 = 10^{-1}, \quad \Gamma_2^0 = 1.$$

### 4.5.2 Numerical discretization

As in [ZCW08b], we approximate equations (4.22)–(4.24) in the one-dimensional domain $\Omega = (0,1)$ by a BDF2 (second-order Backward Differentiation Formula) discretization in time. The spatial discretization is performed by finite volumes. The scheme is explicit for the mobility and potential, using the second-order approximation $\bar{u}^k := 2u^{k-1} - u^{k-2}$, but implicit in the reactions and semi-implicit in the diffusion. Let $\Delta t > 0$ be the time step size, $\Delta x > 0$ the space grid size, and $x_i = i\Delta x$, $x_{i\pm 1/2} = (i\pm 1/2)\Delta x$. We introduce finite-volume cells $K_i = (x_{i-1/2}, x_{i+1/2})$ for $i = 1, \ldots, N_x$. Then the values $u_i^k$, $v_i^k$, and $\mu_i^k$ approximate $u(x_i, k\Delta t)$, $v(x_i, k\Delta t)$, and $\mu(x_i, k\Delta t)$ respectively for $i = 1, \ldots, N_x$, $k = 1, \ldots, N_T$. Our scheme reads for $k \geq 2$ as follows:

$$\frac{\Delta x}{2\Delta t}(3v_i^k - 4v_i^{k-1} + v_i^{k-2}) + \mathcal{G}_{i+1/2}^k - \mathcal{G}_{i-1/2}^k = -\Delta x R_c^0 u_i^k v_i^k,$$

$$\frac{\Delta x}{2\Delta t}(3u_i^k - 4u_i^{k-1} + u_i^{k-2}) + \mathcal{F}_{i+1/2}^k - \mathcal{F}_{i-1/2}^k = \Delta x u_i^k (1 - u_i^k)\frac{R_p^0 v_i^k}{K + v_i^k},$$

$$\mathcal{H}_{i+1/2}^k - \mathcal{H}_{i-1/2}^k + \Delta x f'(\bar{u}_i^k) = \Delta x \mu_i^k,$$

where the numerical fluxes are given by

$$\mathcal{G}_{i+1/2}^k = -D_0(1 - u_{i+1/2}^k)\frac{v_{i+1}^k - v_i^k}{\Delta x},$$

$$\mathcal{F}_{i+1/2}^k = -M_0 u_{i+1/2}^k (1 - u_{i+1/2}^k)\frac{\mu_{i+1}^k - \mu_i^k}{\Delta x}, \quad \mathcal{H}_{i+1/2}^k = -\frac{u_{i+1}^k - u_i^k}{\Delta x},$$

and $u_{i+1/2}^k = \frac{1}{2}(u_{i+1}^k + u_i^k)$. The approximation $(u_i^1, v_i^1, \mu_i^1)$ at the first time step is computed from the implicit Euler method.

In the same way, we discretized a simplified version of [WZ12] which reads in its dimensionless form for $k \geq 2$ as

$$\frac{\Delta x}{2\Delta t}(3w_i^k - 4w_i^{k-1} + w_i^{k-2}) + \mathcal{G}_{i+1/2}^k - \mathcal{G}_{i-1/2}^k = -\Delta x u_i^k \frac{\widetilde{R}_c^0 v_i^k}{\widetilde{K} + v},$$

$$\frac{\Delta x}{2\Delta t}(3u_i^k - 4u_i^{k-1} + u_i^{k-2}) + \widetilde{\mathcal{F}}_{i+1/2}^k - \widetilde{\mathcal{F}}_{i-1/2}^k = \Delta x u_i^k \frac{R_p^0 v_i^k}{K + v_i^k},$$

$$\mathcal{H}_{i+1/2}^k - \mathcal{H}_{i-1/2}^k + \Delta x f'(\bar{u}_i^k) = \Delta x \mu_i^k,$$

where we abbreviated $w_i^k = (1 - u_i^k)v_i^k$, $\mathcal{G}$ and $\mathcal{H}$ are as above, $\widetilde{\mathcal{F}} = -M_0 u_{i+1/2}^k (\mu_{i+1}^k - \mu_i^k)/\Delta x$, and $\widetilde{R}_c^0 = 1$, $R_p^0 = 1$ are scaled rates. We use the Newton method to solve the resulting system of nonlinear equations. For the first three test cases, we used a mesh of 128 cells and the time step size $\Delta t = 10^{-3}$.

### 4.5.3 Numerical results

**Test case** 1**:**

We consider the initial conditions

$$u^0(x) = \frac{1}{2}\sin(2\pi x)^2 + 2 \cdot 10^{-2}, \quad v^0(x) \equiv 0.75.$$

The numerical solutions $u$ and $v$ are presented in Figure 4.1. The substrate concentration converges uniformly to zero as $t \to \infty$ because of the consumption term, while the volume fraction of the biomass is increasing in time. The increase becomes slower and stops after some time since the production term is proportional to the substrate concentration which almost vanishes for large times and hence the production term vanishes too. In our model, both the biomass fraction and the substrate concentration change at a slower rate compared to the model of [WZ12], which is caused by the additional factor $1 - u$ in the source term. Accordingly, the convergence to the steady state is smaller in our model than in the model of [WZ12]. Note that, without the additional factor $1 - u$, an initial value $u^0$ smaller but close to one may lead to a volume fraction exceeding its maximal value and consequently break down the numerical scheme.
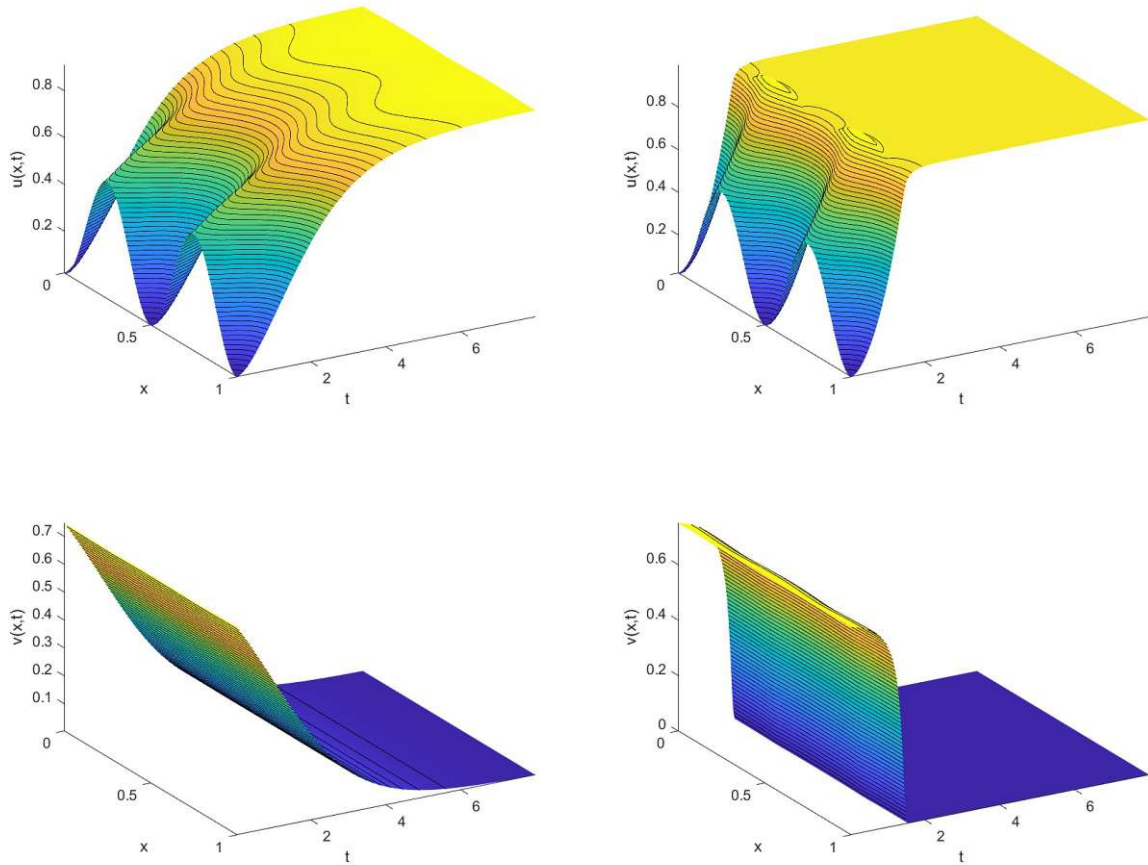
Figure 4.1: Biomass fraction $u$ (top) and substrate concentration $v$ (bottom) in test case 1 for our system (left) and the system of [WZ12] (right).

**Test case** 2:

We consider the initial conditions

$$u^0(x) = \begin{cases} 0.2, & \text{if } 0 \le x \le 0.2, \\ 1 \cdot 10^{-2}, & \text{if } 0.2 < x \le 1, \end{cases} \quad v^0(x) \equiv 0.1.$$

In both models, the volume fraction of biomass growths rather fast until the substrate concentration vanishes; see Figure 4.2. Due to the additional factor $1 - u$ in our mobility, we can observe a slower diffusion in areas of larger volume fraction compared to [WZ12]. In areas of low volume fractions, we observe a larger growth than for [WZ12], which can be explained by the larger nutrient consumption compared to our model, causing a lack of nutrient supply for further growth.

Figure 4.2: Biomass $u$ in test case 2 for our system (left) and the system of [WZ12] (right).

**Test case** $3$**:**

We choose the initial conditions

$$u^0(x) = -(x - 1/2)^2 + 1/3, \quad v^0(x) \equiv 0.3. \tag{4.25}$$

As in the previous test cases, we observe in Figure 4.3 a faster growth of biomass volume fraction in the model of [WZ12]. Moreover, the growth process dominates before the diffusion process flattens the maximal volume fraction towards the steady state. Due to the absence of the factor $1 - u$, this effect is stronger than in the model of [WZ12].
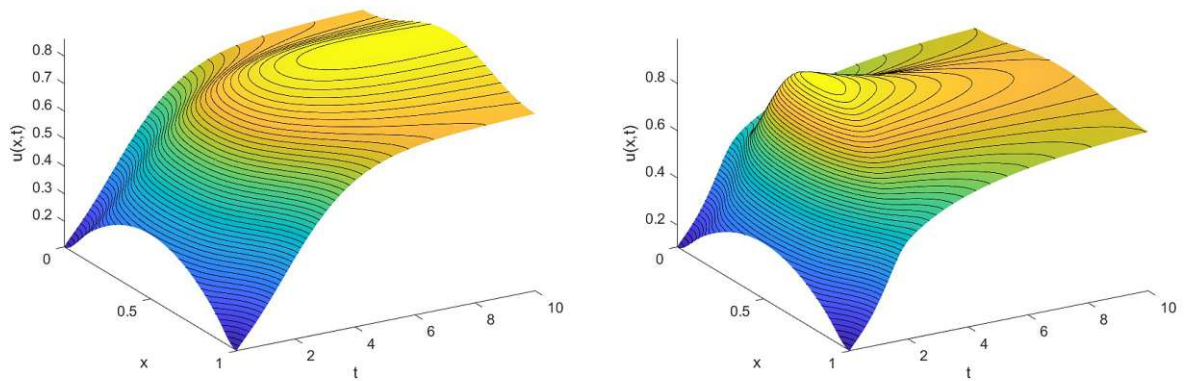


Figure 4.3: Biomass $u$ in test case 3 for our system (left) and the system of [WZ12] (right).

**Test case** 4**:**

We analyze the order of convergence in space with the initial conditions (4.25). Since there does not exist an explicit solution, we compute a reference solution $(u_{\text{ref}}, v_{\text{ref}})$ at time $T = 1$ on a mesh with 2048 cells with time step size $\Delta t = 10^{-5}$. The approximate solutions $u^{(j)}$ are determined on meshes of $2^j$ cells for $j = 4, \ldots, 10$. We choose a rather small value for $T$ to compute the order of convergence in space before a steady state is reached. Figure 4.4 (left) illustrates the discrete $L^2$ norm of the difference $u_{\text{ref}} - u^{(j)}$ for $j = 4, \ldots, 10$. As expected, we observe a second-order convergence in space.



Figure 4.4: Convergence in space (left) and convergence in time (right) at time $T = 1$.

**Test case** 5**:**

We analyze the order of convergence in time by using as before the initial conditions (4.25) and by choosing $L = 128$ cells in space. We compute a reference solution $(u_{\text{ref}}, v_{\text{ref}})$ at time $T = 1$ with time step size $\Delta t = 1/(2^{14}L) \approx 5 \cdot 10^{-7}$. The approximate solutions $u^{(j)}$ are determined with time step sizes $\Delta t = 1/(2^{2j}L)$ for $j = 1, \ldots, 6$. Figure 4.4 (right) illustrates the discrete $L^2$ norm of the difference $u_{\text{ref}} - u^{(j)}$ for $j = 1, \ldots, 6$. We observe a convergence in time of order 1.73 for $u$ and 2 for $v$, respectively.
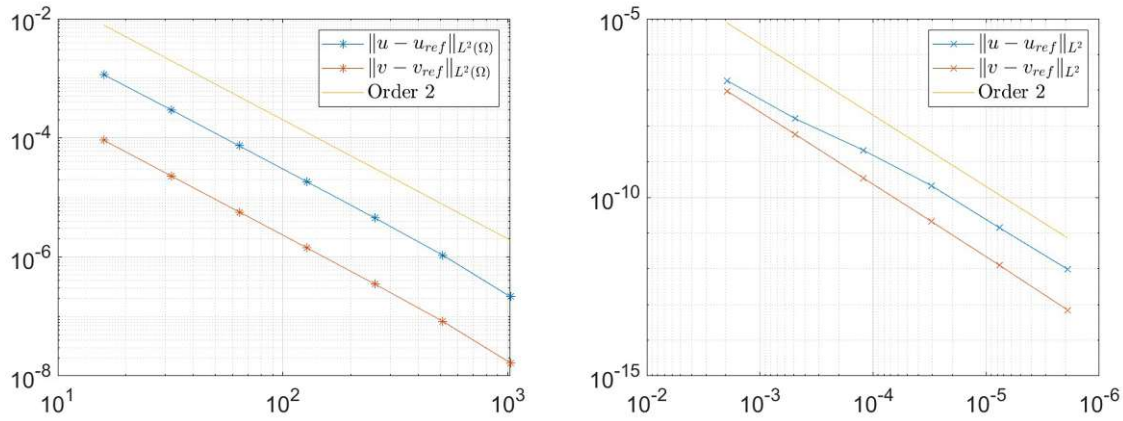
# 5 Discussion and Outlook

We briefly give a discussion of our results and an outlook over possible extensions and open problems connected to this thesis.

## 5.1 The Maxwell–Stefan–Fourier System

We have proved the global existence of weak solutions for the Maxwell–Stefan Fourier system by adapting the techniques of [Jün15]. As pointed out in [JG23], the modeling in our model was inaccurate from a thermodynamical point of view: First and foremost, to be consistent with the Onsager reciprocity principle, the Onsager Matrix should be symmetric. This is considered for the coefficients of the mobility matrix $(M_{ij})_{i,j=1,\ldots,n}$, but not for the terms describing the Soret/Dufour effects. Furthermore, considering a vanishing barycentric velocity, the pressure needs to be considered to be consistent with thermodynamics. Then, as shown in [JG23], the Onsager matrix is positive definite on the space $L = \{\boldsymbol{y} \in \mathbb{R}^n \colon \boldsymbol{y} \cdot \sqrt{\boldsymbol{\rho}}\}$, while we assumed positive definiteness on $L = \{\boldsymbol{y} \in \mathbb{R}^n \colon \boldsymbol{y} \cdot \sqrt{\boldsymbol{1}}\}$. However, from a mathematical point of view, these are only minor shortcomings, as the improved thermodynamical modeling heavily relies on the same estimates as in chapter 2.

Of course, there is room for further research. For instance, the question of existence in case of nonvanishing barycentric velocity is open. Furthermore, the equivalence between the Maxwell–Stefan and the Fick–Onsager formulation remains open in the nonisothermal setting. In Proposition 5, we showed based on [BD23] that the Maxwell–Stefan formulation leads to the Fick–Onsager formulation for a specific choice of coefficients, but we were not able to establish equivalence. Another question which remains open is the longterm behavior of the system. From a practical point of view, numerical analysis for the (improved) model, including simulations and comparisons to experiments would be interesting.

## 5.2 The Quorum Sensing Biofilm Model

We showed the existence, uniqueness and convergence of discrete solutions for the quorum sensing biofilm model of [EHKE15]. The discrete solutions preserve the upper bound for the biomass fraction and the solvent, and furthermore guarantees upper bounds for the autoinducer molecule and the dispersed cells. According to our numerical simulations however, the upper bounds for autoinducer molecules and signal cells overestimate the actual concentration. Hence, it might be possible to improve the upper bounds significantly. Another question worth investigating regards the regularity of the limit. We have proved convergence towards a weak solution $(M, N, S, A)$, yet according to [ESE17], the uniqueness only holds for a smaller class of functions and we can not conclude the convergence of

the whole sequence. Therefore, a proof of the required regularity for the limit $(M, N, S, A)$ would be interesting. As the numerical results showed, we find for the biomass fraction an order of convergence in space of $\approx 1.7$ in the $L^2$–Norm, instead of an expected second order convergence. This might be caused by the degenerate–singular diffusion term, but it would be interesting to analytically prove which order of convergence should actually be obtained. Since biofilms are highly complex, more effects could be taken into account. For instance an extension of the model which includes treatment of the biofilm with biocide could be considered and the numerical analysis could be adapted. Besides, the initial biofilm growth model of [EPL01] is coupled to an incompressible Navier–Stokes equation. This coupling is neglected in the existence analysis [EZE09], as well as in the numerical analysis/simulations. In the existence analysis, this coupling complicates the problem tremendously, as the Navier–Stokes equation is assumed to hold in the bulk liquid region, i.e. $\Omega_b := \{x \in \Omega \mid M(t,x) = 0\}$. Consequently, the sharp biomass front operates as free boundary for the incompressible Navier–Stokes equation. It is yet unclear how to prove the existence of solutions in this case. From a modeling point of view, it would be interesting to compare the biofilm growth model of [EPL01] with the growth model discussed in chapter 4.

## 5.3 Cahn–Hilliard Type Biofilm Growth Model

We discussed a biofilm growth model motivated by [WZ12] and proved the global existence of weak solutions for a modified version of this model by truncating the equations and using a Galerkin approximation.
From a model point of view however, this can only be seen as a first step towards an existence analysis for the binary fluid biofilm models, as the modifications which we described in section 1.4.1 went quite far from the original model. This raises plenty of questions for further research. The model of [WZ12] used the mobility $M(u) = u$. From an analytical point of view however, it does not seem possible to treat mobilities which do not cancel the singularities in $M(u)f''(u)$. This is caused by the degeneracy of the mobility, which requires additional truncations/regularizations in the analysis and the lack of estimates to identify the limit in the deregularization.
Another step back towards the original model of [WZ12] would be, to find a way to treat the additional factor $(1-u)$ in the time derivative of the substrate equation. However, new techniques might be needed as this adds another degeneracy to the equation. The addition of elastic effects, is another possible step towards the original model [WZ12]. Having said that, it complicates the analysis to a great deal as it adds an additional coupling to an Smoluchowski equation (see [WZ12, Equation (16)]) and it yet remains unclear how to prove the global existence.

Apart from the connection to biofilm models, another interesting question is the relation between local and nonlocal Cahn–Hilliard equations. As described in section 1.4.3, some papers investigated the connections between local and nonlocal Cahn–Hilliard equations. This poses the question, whether we could obtain the solution of the model (1.24)–(1.29) as a limit of a nonlocal counterpart.
As the existence analysis for the model of [WZ12] appears to be highly nontrivial without

modifications and tremendous simplifications, this also raises the question whether we can find a potentially more complex yet solvable (from an analytical point of view) nonlocal biofilm growth model.

# Bibliography

[ABSS20]  B. Anwasia, M. Bisi, F. Salvarani, and A. J. Soares. On the Maxwell-Stefan diffusion limit for a reactive mixture of polyatomic gases in non-isothermal setting. *Kinet. Relat. Models*, 13(1):63–95, 2020.

[ACM17]  Boris Andreianov, Clément Cancès, and Ayman Moussa. A nonlinear time compactness result and applications to discretization of degenerate parabolic-elliptic PDEs. *J. Funct. Anal.*, 273(12):3633–3670, 2017.

[AES18]  Md. Afsar Ali, Hermann J. Eberl, and Rangarajan Sudarsan. Numerical solution of a degenerate, diffusion reaction based biofilm growth model on structured non-orthogonal grids. *Commun. Comput. Phys.*, 24(3):695–741, 2018.

[AKK11]  Dimitra C. Antonopoulou, Georgia D. Karali, and Georgios T. Kossioris. Asymptotics for a generalized Cahn-Hilliard equation with forcing terms. *Discrete Contin. Dyn. Syst.*, 30(4):1037–1054, 2011.

[AL83]  Hans Wilhelm Alt and Stephan Luckhaus. Quasilinear elliptic-parabolic differential equations. *Math. Z.*, 183(3):311–341, 1983.

[BB21]  Andrea Bondesan and Marc Briant. Stability of the Maxwell-Stefan system in the diffusion asymptotics of the Boltzmann multi-species equation. *Comm. Math. Phys.*, 382(1):381–440, 2021.

[BCCHF15]  Marianne Bessemoulin-Chatard, Claire Chainais-Hillairet, and Francis Filbet. On discrete functional inequalities for some finite volume schemes. *IMA J. Numer. Anal.*, 35(3):1125–1149, 2015.

[BCH13]  Konstantin Brenner, Clément Cancès, and Danielle Hilhorst. Finite volume approximation for an immiscible two-phase flow in porous media with discontinuous capillary pressure. *Comput. Geosci.*, 17(3):573–597, 2013.

[BD15]  Dieter Bothe and Wolfgang Dreyer. Continuum thermodynamics of chemically reacting fluid mixtures. *Acta Mech.*, 226(6):1757–1805, 2015.

[BD21]  Dieter Bothe and Pierre-Etienne Druet. Mass transport in multicomponent compressible fluids: local and global well-posedness in classes of strong solutions for general class-one models. *Nonlinear Anal.*, 210:Paper No. 112389, 53, 2021.

[BD23]  Dieter Bothe and Pierre-Étienne Druet. On the structure of continuum thermodynamical diffusion fluxes—a novel closure scheme and its relation to the

Maxwell-Stefan and the Fick-Onsager approach. *Internat. J. Engrg. Sci.*, 184:Paper No. 103818, 33, 2023.

[BG20]    Marc Briant and Bérénice Grec. Rigorous derivation of the fick cross-diffusion system from the multi-species boltzmann equation in the diffusive scaling. *arXiv: Analysis of PDEs*, 2020.

[BGP19]    Laurent Boudin, Bérénice Grec, and Vincent Pavan. Diffusion models for mixtures using a stiff dissipative hyperbolic formalism. *J. Hyperbolic Differ. Equ.*, 16(2):293–312, 2019.

[BGPS13]    Laurent Boudin, Bérénice Grec, Milana Pavić, and Francesco Salvarani. Diffusion asymptotics of a kinetic model for gaseous mixtures. *Kinet. Relat. Models*, 6(1):137–157, 2013.

[BJPZ22]    Miroslav Bulíček, Ansgar Jüngel, Milan Pokorný, and Nicola Zamponi. Existence analysis of a stationary compressible fluid model for heat-conducting and chemically reacting mixtures. *J. Math. Phys.*, 63(5):Paper No. 051501, 48, 2022.

[Bot11]    Dieter Bothe. On the Maxwell-Stefan approach to multicomponent diffusion. In *Parabolic problems*, volume 80 of *Progr. Nonlinear Differential Equations Appl.*, pages 81–93. Birkhäuser/Springer Basel AG, Basel, 2011.

[Bry08]    James D. Bryers. Medical biofilms. *Biotechnology and Bioengineering*, 100(1):1–18, 2008.

[CCHGJ19]    Clément Cancès, Claire Chainais-Hillairet, Anita Gerstenmayer, and Ansgar Jüngel. Finite-volume scheme for a degenerate cross-diffusion model motivated from ion transport. *Numer. Methods Partial Differential Equations*, 35(2):545–575, 2019.

[CES23]    José A Carrillo, Charles Elbar, and Jakub Skrzeczkowski. Degenerate cahn-hilliard systems: From nonlocal to local. *arXiv preprint arXiv:2303.11929*, 2023.

[CHLP03]    Claire Chainais-Hillairet, Jian-Guo Liu, and Yue-Jun Peng. Finite volume scheme for multi-dimensional drift-diffusion equations and convergence analysis. *M2AN Math. Model. Numer. Anal.*, 37(2):319–338, 2003.

[CJ15]    Xiuqing Chen and Ansgar Jüngel. Analysis of an incompressible Navier-Stokes-Maxwell-Stefan system. *Comm. Math. Phys.*, 340(2):471–497, 2015.

[CMZ14]    Laurence Cherfils, Alain Miranville, and Sergey Zelik. On a generalized Cahn-Hilliard equation with biological applications. *Discrete Contin. Dyn. Syst. Ser. B*, 19(7):2013–2026, 2014.

[DD21]    Michele Dolce and Donatella Donatelli. Artificial compressibility method for the Navier-Stokes-Maxwell-Stefan system. *J. Dynam. Differential Equations*, 33(1):35–62, 2021.

[DDGG20]   Wolfgang Dreyer, Pierre-Étienne Druet, Paul Gajewski, and Clemens Guhlke. Analysis of improved Nernst-Planck-Poisson models of compressible isothermal electrolytes. *Z. Angew. Math. Phys.*, 71(4):Paper No. 119, 68, 2020.

[Dei85]   Klaus Deimling. *Nonlinear functional analysis.* Springer-Verlag, Berlin, 1985.

[DJ12]   Michael Dreher and Ansgar Jüngel. Compact families of piecewise constant functions in $L^p(0, T; B)$. *Nonlinear Anal.*, 75(6):3072–3077, 2012.

[DJZ21]   Esther S. Daus, Ansgar Jüngel, and Antoine Zurek. Convergence of a finite-volume scheme for a degenerate-singular cross-diffusion system for biofilms. *IMA J. Numer. Anal.*, 41(2):935–973, 2021.

[DMZ19]   Esther S. Daus, Pina Milišić, and Nicola Zamponi. Analysis of a degenerate and singular volume-filling cross-diffusion system modeling biofilm growth. *SIAM J. Math. Anal.*, 51(4):3569–3605, 2019.

[Dru16]   P.-E. Druet. Analysis of improved nernst–planck–poisson models of isothermal compressible electrolytes subject to chemical reactions: The case of a degenerate mobility matrix. *WIAS Preprint no. 2321*, 2016.

[Dru21]   Pierre-Etienne Druet. A theory of generalised solutions for ideal gas mixtures with Maxwell-Stefan diffusion. *Discrete Contin. Dyn. Syst. Ser. S*, 14(11):4035–4067, 2021.

[Dru22]   Pierre-Etienne Druet. Maximal mixed parabolic–hyperbolic regularity for the full equations of multicomponent fluid dynamics. *Nonlinearity*, 35(7):3812, jun 2022.

[Ebe20]   Matthias Ebenbeck. *Cahn-Hilliard-Brinkman models for tumour growth: Modelling, analysis and optimal control.* PhD thesis, 2020.

[EEWZ14]   Hermann J. Eberl, Messoud A. Efendiev, Dariusz Wrzosek, and Anna Zhigun. Analysis of a degenerate biofilm model with a nutrient taxis term. *Discrete Contin. Dyn. Syst.*, 34(1):99–119, 2014.

[EG96]   Charles M. Elliott and Harald Garcke. On the Cahn-Hilliard equation with degenerate mobility. *SIAM J. Math. Anal.*, 27(2):404–423, 1996.

[EGH00]   Robert Eymard, Thierry Gallouët, and Raphaèle Herbin. Finite volume methods. In *Handbook of numerical analysis, Vol. VII*, Handb. Numer. Anal., VII, pages 713–1020. North-Holland, Amsterdam, 2000.

[EGHM02]   Robert Eymard, Thierry Gallouët, Raphaèle Herbin, and Anthony Michel. Convergence of a finite volume scheme for nonlinear degenerate parabolic equations. *Numer. Math.*, 92(1):41–82, 2002.

[EHKE15]   Blessing O Emerenini, Burkhard A Hense, Christina Kuttler, and Hermann J Eberl. A mathematical model of quorum sensing induced biofilm detachment. *PloS one*, 10(7):e0132385, 2015.

[EPL01]     Herman. J. Eberl, David F. Parker, and Mark C. M. Van Loosdrecht. A new deterministic spatio-temporal continuum model for biofilm development. *Journal of Theoretical Medicine*, 3(3):161–175, 2001.

[ESE17]     Blessing O. Emerenini, Stefanie Sonner, and Hermann J. Eberl. Mathematical analysis of a quorum sensing induced biofilm dispersal model and numerical simulation of hollowing effects. *Math. Biosci. Eng.*, 14(3):625–653, 2017.

[EZE09]     Messoud A. Efendiev, Sergey V. Zelik, and Hermann J. Eberl. Existence and longtime behavior of a biofilm model. *Commun. Pure Appl. Anal.*, 8(2):509–531, 2009.

[FAG09]     Pina M Fratamico, Bassam A Annous, and NW Guenther. *Biofilms in the food and beverage industries*. Elsevier, 2009.

[FHKM22]    Julian Fischer, Katharina Hopf, Michael Kniely, and Alexander Mielke. Global existence analysis of energy-reaction-diffusion systems. *SIAM Journal on Mathematical Analysis*, 54(1):220–267, 2022.

[FHX14]     A. Friedman, B. Hu, and C. Xue. On a multiphase multicomponent model of biofilm growth. *Arch. Ration. Mech. Anal.*, 211:257–300, 2014.

[FLR17]     Sergio Frigeri, Kei Fong Lam, and Elisabetta Rocca. On a diffuse interface model for tumour growth with non-local interactions and degenerate mobilities. *Solvability, Regularity, and Optimal Control of Boundary Value Problems for PDEs: In Honour of Prof. Gianni Gilardi*, pages 217–254, 2017.

[Fri21]     Sergio Frigeri. On a nonlocal Cahn-Hilliard/Navier-Stokes system with degenerate mobility and singular potential for incompressible fluids with different densities. *Ann. Inst. H. Poincaré C Anal. Non Linéaire*, 38(3):647–687, 2021.

[GHHK20]    Martin J. Gander, Laurence Halpern, Florence Hubert, and Stella Krell. Optimized overlapping DDFV Schwarz algorithms. In *Finite volumes for complex applications IX—methods, theoretical aspects, examples—FVCA 9, Bergen, Norway, June 2020*, volume 323 of *Springer Proc. Math. Stat.*, pages 365–373. Springer, Cham, [2020] ©2020.

[GL16]      Harald Garcke and Kei Fong Lam. Global weak solutions and asymptotic limits of a cahn–hilliard–darcy system modelling tumour growth. *AIMS Mathematics*, 1(3):318–360, 2016.

[GM98]      Vincent Giovangigli and Marc Massot. The local Cauchy problem for multicomponent reactive flows in full vibrational non-equilibrium. *Math. Methods Appl. Sci.*, 21(15):1415–1439, 1998.

[GPZ15]     Vincent Giovangigli, Milan Pokorný, and Ewelina Zatorska. On the steady flow of reactive gaseous mixture. *Analysis (Berlin)*, 35(4):319–341, 2015.

[HES22]     Jack M. Hughes, Hermann J. Eberl, and Stefanie Sonner. A mathematical model of discrete attachment to a cellulolytic biofilm using random DEs. *Math. Biosci. Eng.*, 19(7):6582–6619, 2022.

[HJ21]      Christoph Helmer and Ansgar Jüngel. Analysis of Maxwell-Stefan systems for heat conducting fluid mixtures. *Nonlinear Anal. Real World Appl.*, 59:Paper No. 103263, 19, 2021.

[HJ23]      Christoph Helmer and Ansgar Jüngel. Existence analysis for a reaction-diffusion cahn-hilliard-type system with degenerate mobility and singular potential modeling biofilm growth, 2023.

[HJT19]     Xiaokai Huo, Ansgar Jüngel, and Athanasios E. Tzavaras. High-friction limits of Euler flows for multicomponent systems. *Nonlinearity*, 32(8):2875–2913, 2019.

[HJZ23]     Christoph Helmer, Ansgar Jüngel, and Antoine Zurek. Analysis of a finite-volume scheme for a single-species biofilm model. *Appl. Numer. Math.*, 185:386–405, 2023.

[HMPW17]    Martin Herberg, Martin Meyries, Jan Prüss, and Mathias Wilke. Reaction-diffusion systems of Maxwell-Stefan type with reversible mass-action kinetics. *Nonlinear Anal.*, 159:264–284, 2017.

[HS18]      Harsha Hutridurga and Francesco Salvarani. Existence and uniqueness analysis of a non-isothermal cross-diffusion system of Maxwell-Stefan type. *Appl. Math. Lett.*, 75:108–113, 2018.

[HSCS04]    Luanne Hall-Stoodley, J William Costerton, and Paul Stoodley. Bacterial biofilms: from the natural environment to infectious diseases. *Nature reviews microbiology*, 2(2):95–108, 2004.

[IM18]      Annalisa Iuorio and Stefano Melchionna. Long-time behavior of a non-local Cahn-Hilliard equation with reaction. *Discrete Contin. Dyn. Syst.*, 38(8):3765–3788, 2018.

[JG23]      Ansgar Jüngel and Stefanos Georgiadis. Global Existence of Weak Solutions and Weak–Strong Uniqueness for Nonisothermal Maxwell–Stefan Systems. 2023.

[JL19]      Ansgar Jüngel and Oliver Leingang. Convergence of an implicit Euler Galerkin scheme for Poisson-Maxwell-Stefan systems. *Adv. Comput. Math.*, 45(3):1469–1498, 2019.

[JS13]      Ansgar Jüngel and Ines Viktoria Stelzer. Existence analysis of Maxwell-Stefan systems for multicomponent mixtures. *SIAM J. Math. Anal.*, 45(4):2421–2440, 2013.

[Jün15]     Ansgar Jüngel. The boundedness-by-entropy method for cross-diffusion systems. *Nonlinearity*, 28(6):1963–2001, 2015.

[JZ21]     Ansgar Jüngel and Antoine Zurek. A convergent structure-preserving finite-volume scheme for the shigesada–kawasaki–teramoto population system. *SIAM Journal on Numerical Analysis*, 59(4):2286–2309, 2021.

[JZ22]     Ansgar Jüngel and Antoine Zurek. A discrete boundedness-by-entropy method for finite-volume approximations of cross-diffusion systems. *IMA Journal of Numerical Analysis*, 43(1):560–589, 01 2022.

[KD10]     Isaac Klapper and Jack Dockery. Mathematical description of microbial biofilms. *SIAM Review*, 52(2):221–265, 2010.

[Lud56]    Carl Ludwig. Diffusion zwischen ungleich erwärmten orten gleich zusammengesetzer lösungen. *Sitzungsberichte der Mathematisch-Naturwissenschaftlichen Classe der Kaiserlichen Akademie der Wissenschaften Wien, 2te Abteilung*, 20, 1856.

[Max67]    J. Clerk Maxwell. On the dynamical theory of gases. *Philosophical Transactions of the Royal Society of London*, 157:49–88, 1867.

[ME80]     Robert G Mortimer and Henry Eyring. Elementary transition state theory of the soret and dufour effects. *Proceedings of the National Academy of Sciences of the United States of America*, 77(4):1728–1731, 1980.

[MR15]     Stefano Melchionna and Elisabetta Rocca. On a nonlocal cahn-hilliard equation with a reaction term. *arXiv preprint arXiv:1501.01541*, 2015.

[MRST19]   Stefano Melchionna, Helene Ranetbauer, Luca Scarpa, and Lara Trussardi. From nonlocal to local Cahn-Hilliard equation. *Adv. Math. Sci. Appl.*, 28(2):197–211, 2019.

[NC08]     Amy Novick-Cohen. Chapter 4 the cahn–hilliard equation. volume 4 of *Handbook of Differential Equations: Evolutionary Equations*, pages 201–228. North-Holland, 2008.

[Ons31]    Lars Onsager. Reciprocal relations in irreversible processes. i. *Phys. Rev.*, 37:405–426, Feb 1931.

[OR20]     Lukas Ostrowski and Christian Rohde. Compressible multicomponent flow in porous media with Maxwell-Stefan diffusion. *Math. Methods Appl. Sci.*, 43(7):4200–4221, 2020.

[PMCW11]   Steven L Percival, Sladjana Malic, Helena Cruz, and David W Williams. Introduction to biofilms. *Biofilms and veterinary medicine*, pages 41–68, 2011.

[Poz18]    Jose Luis Del Pozo. Biofilm-related disease. *Expert Review of Anti-infective Therapy*, 16(1):51–65, 2018. PMID: 29235402.

[PP17]     Tomasz Piasecki and Milan Pokorný. Weak and variational entropy solutions to the system describing steady flow of a compressible reactive mixture. *Nonlinear Anal.*, 159:365–392, 2017.

[PP21]       Benoît Perthame and Alexandre Poulain. Relaxation of the Cahn-Hilliard equation with singular single-well potential and degenerate mobility. *European J. Appl. Math.*, 32(1):89–112, 2021.

[PVLH98a]   Cristian Picioreanu, Mark CM Van Loosdrecht, and Joseph J Heijnen. Mathematical modeling of biofilm structure with a hybrid differential-discrete cellular automaton approach. *Biotechnology and bioengineering*, 58(1):101–116, 1998.

[PvLH98b]   Cristian Picioreanu, Mark CM van Loosdrecht, and Joseph J Heijnen. A new combined differential-discrete cellular automaton approach for biofilm modeling: Application for growth in gel beads. *Biotechnology and bioengineering*, 57(6):718–731, 1998.

[PvLH99]    Cristian Picioreanu, M van Loosdrecht, and J Heijnen. *Multidimensional modeling of biofilm structure*. Delft University of Technology, Faculty of Applied Sciences, 1999.

[RM80]      Bruce E. Rittmann and Perry L. McCarty. Evaluation of steady-state-biofilm kinetics. *Biotechnology and Bioengineering*, 22(11):2359–2373, 1980.

[RSE15]     Kazi A. Rahman, Rangarajan Sudarsan, and Hermann J. Eberl. A mixed-culture biofilm model with cross-diffusion. *Bull. Math. Biol.*, 77(11):2086–2124, 2015.

[SEL14]     Cristina Solano, Maite Echeverz, and Iñigo Lasa. Biofilm dispersion and quorum sensing. *Current Opinion in Microbiology*, 18:96–104, 2014. Cell regulation.

[Sim87]     Jacques Simon. Compact sets in the space $L^p(0, T; B)$. *Ann. Mat. Pura Appl. (4)*, 146:65–96, 1987.

[SN16]      Shama Sehar and Iffat Naz. Role of the biofilms in wastewater treatment. *Microbial biofilms-importance and applications*, pages 121–144, 2016.

[Ste71]     Josef Stefan. Über das gleichgewicht und die bewegung, insbesondere die diffusion von gasgemengen. *Sitzungsberichte der Mathematisch-Naturwissenschaftlichen Classe der Kaiserlichen Akademie der Wissenschaften Wien, 2te Abteilung*, 63:63 – 124, 1871.

[TA99]      Shigeru Takata and Kazuo Aoki. Two-surface problems of a multicomponent mixture of vapors and noncondensable gases in the continuum limit in the light of kinetic theory. *Physics of Fluids*, 11(9):2743–2756, 1999.

[Tem97]     Roger Temam. *Infinite-dimensional dynamical systems in mechanics and physics*, volume 68 of *Applied Mathematical Sciences*. Springer-Verlag, New York, second edition, 1997.

[TK93]      R. Taylor and R. Krishna. *Multicomponent Mass Transfer*. Wiley Series in Chemical Engineering. Wiley, 1993.

[WG86]     O. Wanner and W. Gujer. A multispecies biofilm model. *Biotechnology and Bioengineering*, 28(3):314–328, 1986.

[WZ10]     Qi Wang and Tianyu Zhang. Review of mathematical models for biofilms. *Solid State Communications*, 150(21):1009–1022, 2010. Nanoscale Interfacial Phenomena in Complex Fluids.

[WZ12]     Qi Wang and Tianyu Zhang. Kinetic theories for biofilms. *Discrete Contin. Dyn. Syst. Ser. B*, 17(3):1027–1059, 2012.

[Yin92]    Jing Xue Yin. On the existence of nonnegative continuous solutions of the Cahn-Hilliard equation. *J. Differential Equations*, 97(2):310–327, 1992.

[ZCW08a]   Tianyu Zhang, N. G. Cogan, and Qi Wang. Phase field models for biofilms. I. Theory and one-dimensional simulations. *SIAM J. Appl. Math.*, 69(3):641–669, 2008.

[ZCW08b]   Tianyu Zhang, N. G. Cogan, and Qi Wang. Phase-field models for biofilms ii 2-d numerical simulations of biofilm-flow interaction. *Communications in Computational Physics*, 4(1):72–101, 2008.