

Some Implications of Constraints in Phasefield Models

Zur Erlangung des akademischen Grades eines
Doktors der Ingenieurwissenschaften (Dr.-Ing.)
von der KIT-Fakultät für Maschinenbau des
Karlsruher Instituts für Technologie (KIT)

genehmigte
Dissertation
von

Dipl.-Ing. Andreas Reiter

Tag der mündlichen Prüfung: 04.12.2020
Erstgutachterin: Prof. Dr. rer. nat. Britta Nestler
Zweitgutachterin: Prof. Dr. Katharina Schratz

Kurzfassung

In dieser Arbeit werden verschiedene Aspekte, die sich durch Zwangsbedingungen in der Phasenfeldmodellierung ergeben, untersucht.

Zum einen wird, im Rahmen eines reinen Phasenfeldmodells, der Einfluss des häufig verwendeten Hindernispotentials in Bezug auf die Diskretisierung und algorithmische Gesichtspunkte der Verwendung von Projektions-basierten Algorithmen in nicht-gewichteten und gewichteten Mobilitätsformulierungen betrachtet.

Zum anderen werden "Grandchem"-artige Modelle in einem chemischen, mechanischen und chemo-mechanischem Kontext diskutiert, in denen eine gegebene phasenunabhängige Größe innerhalb von Mehrphasenbereichen als gewichtetes Mittel der entsprechenden Größen innerhalb der Einzelphasen aufgefasst wird. Die so eingeführten zusätzlichen Freiheitsgrade ermöglichen durch eine geschickte Festlegung der phasenspezifischen Werte in Abhängigkeit der restlichen Parameter eine verbesserte Modellbildung, durch welche sich der Einfluss der Breite der Übergangsbereiche auf die Ergebnisse deutlich reduzieren lässt. In vielen Fällen lässt sich die meistens direkt physikalisch motivierte Festlegung der phasenspezifischen Größen zugleich als die Lösung eines parametrisierten Minimierungs- oder Maximierungsproblems unter der Nebenbedingung des vorgegebenen Mittelwerts interpretieren. Hier wird untersucht, welche Konsequenzen sich aus dieser Interpretation ergeben und weshalb das Zusammenspiel dieses lokalen Extremalproblems mit dem globalen variationellen Ansatz des Phasenfeldmodells von entscheidender Bedeutung ist.

Abstract

In this work, several aspects resulting from constraints in phasefield models are investigated.

On the one hand, within the framework of a pure phasefield model, the influence of the commonly used obstacle potential is considered, both with respect to the discretization and to some algorithmic implications for projection-based algorithms using nonweighted and weighted mobility formulations.

On the other hand, “grandchem”-type models, which are based upon interpreting a given phase-independent quantity within the multiphase regions as a weighted average of the corresponding quantity within the individual phases, are discussed within a chemical, mechanical and chemomechanical context. The additional degrees of freedom introduced by these phasespecific quantities lead, if fixed in an clever manner in terms of the remaining parameters, to improved models with a significantly reduced dependence upon the width of the diffuse interface. In many cases, this typically physically motivated specification of the phasespecific quantities can also be interpreted as the solution of a parameterized minimization or maximization problem under the constraint of maintaining the correct average. Here, the investigation focuses on consequences of this interpretation and the importance of the interplay of this local extremum problem with the global variational Ansatz of the phasefield model.

Contents

I	Introduction and Outline	3
1	Introduction	4
2	A Quick Sketch of the Phasefield Method	6
3	An Outline of the Topics Considered	9
3.1	Pure Phasefield Problems	9
3.2	Coupled Problems	10
3.2.1	Solidification Problems	11
3.2.2	Solid-Solid Phase Transformation Problems	14
II	Background	16
4	Equality and Inequality Constrained Problems	18
4.1	Linear and Nonlinear Equality Constraints - the Finite-Dimensional Case	18
4.2	Linear and Nonlinear Equality Constraints - the Infinite-Dimensional Case	23
4.3	Linear and Nonlinear Inequality Constraints	27
5	Lagrange Functions and Legendre-Fenchel Duality	32
5.1	Lagrange Functions	32
5.2	The Legendre-Transform	37
III	Applications	44
6	Pure Phasefield Problems	47
6.1	The Basic Phasefield Functional	48
6.2	Two-Phase Problems	53
6.2.1	The Steady-State and Dynamic Phasefield Equations	54
6.2.2	The One-Dimensional Case and the “Standard” Phasefield Profiles	59
6.2.3	Some Analytics for the Discrete Case	65
6.3	Multi-Phasefield Problems	82
6.3.1	The Steady-State Equations	82
6.3.2	The Choice of Dynamics	86
6.3.3	Some Numerical and Algorithmic Considerations	88
7	Applications in the Material Sciences	111
7.1	Quantitative Phasefield Models for Solidification	113
7.1.1	The “Traditional” Model	113
7.1.2	The “More Advanced” Free Energy Model in the Isothermal Case	114

7.1.3	The Resulting Driving Force	120
7.1.4	The Relation with the Grand Chemical Potential	121
7.1.5	The Practical Evaluation of $f(\phi, \mathbf{c}, T)$ and the Chemical Potential	130
7.1.6	An Extension to the Non-Isothermal Case	139
7.2	Quantitative Models for Solid-Solid Transformations	149
7.2.1	“Traditional” Phasefield Models	152
7.2.2	The Quantitative Model in the Two-Phase Case	158
7.2.3	Comparison of the Different Formulations in the Two-Phase Case	164
7.2.4	Some Partial Extension to the Multiphase Case	190
7.2.5	Chemo-Elasticity	206
8	Summary	216
A	Calculation of the Discrete One-Dimensional Phasefield Energy and Interface Width	218
A.1	A Quick Recap of the First-Order Analysis	218
A.2	The Discrete Energetic Analysis	219
A.3	The Local Analysis and the Second-Order Conditions	225
	Bibliography	230

Part I

Introduction and Outline

Chapter 1

Introduction

Over the course of the last decades, the phasefield method has evolved from a popular model in material science problems, in particular for solidification problems in which surface tensions play a major role, to a more widely used computational method for a large class of multiphysics interface and free discontinuity problems. Its primary advantage in problems involving moving interfaces is that, whereas sharp interface approaches typically require some relatively expensive interface tracking, the phasefield method is in many cases capable of reproducing the essential physical effects associated with the interfaces without requiring an explicit identification of their location. This is primarily achieved by approximating the sharp interface with a diffuse transition region with a small but finite width. Within this region, the underlying physics are then approximated through appropriate interpolation procedures, which, if done properly, leads to an approximating the solution of the original sharp interface problem. As the interpolation procedures do not require explicitly identifying the precise location of any surfaces, the phasefield method is often referred to as an interface capturing rather than an interface tracking approach.

There is a price to be paid though, as the diffuse interface approximation introduces an additional lengthscale in the form of the width of the interfaces. In order to provide accurate approximations to the solution of sharp interface problems, it is usually necessary to choose an interface width which is significantly smaller than the other lengthscales of the problem. This in turn requires a discretization which is, at least locally, sufficiently fine to faithfully solve the underlying equations and can thus give rise to relatively large and potentially expensive problems. Whether or not the phasefield approach is competitive is thus clearly a question of both how well the physics are approximated in the diffuse setting and of the efficiency with which the numerical problems can be solved.

Two major developments in this context have been the increasingly widespread use of obstacle-based potentials and an improved modeling approach for the physics within the transition region. In contrast to well-based potentials, obstacle-based potentials in principle allow to significantly decrease the computational cost as they ensure that the computationally expensive interfaces are strictly contained within a small region of width $\mathcal{O}(\epsilon)$. Their only disadvantage is that this requires introducing inequality constraints on the phasefield values, which leads to some additional complications in both the equations and various questions related to their solution. Some primarily practical aspects related to this will be discussed in Chapter 6.

A somewhat similar observation can be made for the improved modeling approach first proposed in the chemically driven context by [42]. Its basic idea is that, considering the interface region as a mixture of different phases, the values of any further physically relevant quantity such as e.g. the concentration, temperature or mechanical strains, can also be considered as the average of the respective quantity over all phases present at this point, i.e. to introduce phase-specific versions of these quantities which need to average to the total ones. As there are of course many ways to redistribute a given quantity onto several phase-specific ones, their precise specification

introduces an additional degree of freedom into the model, which can be used to include more of the underlying physics. Even though this is an intuitively appealing idea, it is also clear that this introduces additional complexities into the model as there is now a “submodel” for deriving the phase-specific quantities from the remaining parameters of the problem. Besides a higher computational difficulty of such models, the additional unknowns further leads to some questions regarding the interpretation and consistency with respect to the standard variational approach underlying the phasefield method. These questions, in particular in relation with various models extending the one in [42] to a mechanical setting, will be the focus of Chapter 7.

Chapter 2

A Quick Sketch of the Phasefield Method

In its most basic form, the phasefield method provides a means of approximating the surface, or, more generally, a surface energy associated with subsets of \mathbf{R}^n . If Ω is such a (sufficiently smooth) subset, its surface area is given by $|\partial\Omega| = \int_{\partial\Omega} d\mathbf{s}$ with $d\mathbf{s}$ denoting the surface measure on $\partial\Omega$. In addition, if one assumes that to the surface one can associate a surface energy density γ , $[\gamma] = \frac{\text{J}}{\text{m}^2}$, this surface can, in the simplest case, be associated with a surface energy of $\mathcal{E}_s(\Omega) = \gamma|\partial\Omega|$. More generally, γ may depend on \mathbf{x} and other parameters such as the temperature, leading to the expression

$$\mathcal{E}(\Omega) = \int_{\partial\Omega} \gamma(\mathbf{x}, \dots) d\mathbf{s}. \quad (2.1)$$

An alternative formulation for the surface energy (2.1) may be obtained by introducing the characteristic function

$$\chi_\Omega(\mathbf{x}) = \begin{cases} 1 & , \mathbf{x} \in \Omega, \\ 0 & , \text{else} \end{cases}$$

of Ω . Even though χ is discontinuous on $\partial\Omega$, one can define its distributional derivative $D\chi_\Omega$, whose total variation can be identified with the surface measure on $\partial\Omega$. Formally, one thus has $|\partial\Omega| = \int_{\mathbf{R}^n} |D\chi_\Omega|$ for the surface area, and, by again using the surface energy density γ ,

$$\mathcal{E}(\Omega) = \int \gamma(\mathbf{x}, \dots) |D\chi_\Omega|. \quad (2.2)$$

While the rigorous form of this this approach in terms of functions of bounded variation turns out to be very fruitful from a theoretical point of view (see e.g. [2] and [6] for an introduction to this topic), it is clear that the use of Equation (2.2) in numerical approximations is somewhat difficult due to the lack of smoothness.

A first remedy would be to replace the discontinuous indicator function by a smooth approximation obtained e.g. by convolution with an appropriate smoothing kernel, $\chi_\Omega^\epsilon = \psi_\epsilon * \chi_\Omega$. As the approximation χ_Ω^ϵ is smooth, its distributional derivative reduces to its classical gradient, and thus (2.2) could be approximated by

$$\mathcal{E}_\epsilon(\Omega) = \int \gamma(\mathbf{x}, \dots) |\nabla \chi_\Omega^\epsilon| d\mathbf{x}. \quad (2.3)$$

While this provides a useful approximation of the surface energy **given** an appropriate subset Ω , it is still not particularly convenient for the purpose of actually finding subsets with minimal surfaces. On the one hand, it is not obvious how one can integrate the “constraint” of corresponding

to a regularized characteristic function into the problem of minimizing the approximated surface energy (2.3). On the other hand, the integrand in (2.3) is still not differentiable due to the kink in the mapping $\zeta \mapsto |\zeta|$.

The solution to this dilemma in the phasefield approach is also based on using an approximation ϕ to an indicator function, but by approximating the surface energy through a combination of two competing terms. The first contribution is given by the **gradient energy density** $\epsilon a(\nabla\phi)$ acting as a (smooth) penalization of gradients. The second contribution is by the **bulk potential** $\frac{1}{\epsilon}w(\phi)$ penalizing values different from 0 and 1. By a proper choice of a and w , the approximation to the minimal surface energy problem in the phasefield setting can then be put as that of minimizing the energy functional

$$\mathcal{E}_\epsilon(\phi) := \int_{\Omega} \epsilon a(\nabla\phi) + \frac{1}{\epsilon}w(\phi) \, d\mathbf{x} \quad (2.4)$$

over an appropriate set of admissible functions.

It is clear that any minimizer of the w -term alone would take only the values 0 and 1, and that this term thus forces ϕ towards some indicator function. In contrast, the contribution from the a -term enforces some spatial regularity of ϕ , and has the purpose of “smoothing” the minimizers that would result from w alone, leading to the desired smooth approximation of χ_{Ω} . The different scalings in ϵ in Equation (2.4) provide a means of regulating the balance between the two terms, with decreasing values of ϵ increasingly favoring the w -contribution, therefore leading to a stronger similarity of the minimizers with actual characteristic functions (i.e. smaller transition regions between the two “phases”). As ϵ tends to zero, this approach can further be shown to converge to the desired sharp-interface limit in the appropriate sense¹.

As one now disposes of a formulation for (approximately) describing minimization problems involving surface energies, it is natural to extend the approach to more general problems involving additional competing energy contributions as these can easily be integrated by modifying the energy functional (2.4). A simple such example is obtained by considering melting/solidification processes involving a liquid and a solid phase where, depending on the temperature, it is energetically favorable for the substance to be either in its solid or liquid form. This leads, besides the surface energy associated with liquid-solid interfaces, to an additional energy contribution. By identifying for example ϕ with an approximation of the indicator function for the liquid region, this effect can be included into the approach above by considering the modified energy functional

$$\tilde{\mathcal{E}}_\epsilon(\phi, T) := \int_{\Omega} \epsilon a(\nabla\phi) + \frac{1}{\epsilon}w(\phi) + f(\phi, T) \, d\mathbf{x}. \quad (2.5)$$

Here T is the temperature (for the moment assumed to be **fixed**) and f represents a suitable approximation of the energy associated with being liquid instead of solid for a given temperature, for example expressed in the form $\Delta f^{ls}(T)h(\phi)$, where $\Delta f^{ls}(T)$ is the energy required for changing from the solid to the liquid state at the given temperature and h is a suitable interpolation function.

It is clear that similar ideas can in principle be applied to a relatively large class of problems². The primary use of the phasefield method within the material sciences is for the description

¹This is usually done in terms of de Giorgi’s Γ -convergence, see [50] and [45] for some early works in this direction within the phasefield setting and e.g. [15] or [47] for a general introduction to Γ -convergence. That this limit procedure requires some slightly more complex definition of convergence is of course to be expected, since the desired limit is not part of the original function space and the functional (more specifically the quadratic gradient energy density) is ill-defined on this limit.

²If any additional parameters such as T in the functionals are not fixed, the situation becomes somewhat more delicate, a point which will be discussed in more detail in Section 3.2.

of dynamic processes involving phase-transformations and/or moving interfaces. In these, the process itself is typically considered to play an important role, and one can not content oneself with a “steady-state” analysis based solely on any actual minimizers of some functionals such as in (2.4) or (2.5). While an appropriately chosen thermodynamic potential (e.g. the entropy or free energy) can serve as a guideline and motivation for the description, these potentials are inherently related to equilibrium thermodynamics, i.e. purely spatial. Therefore, by themselves, they do not provide any direct information on the dynamic behavior of the system. Nevertheless, the standard approach within the phasefield setting is a variationally motivated one, i.e. one chooses a suitable thermodynamic potential and then postulates that the dynamics of the independent variables are proportional to the gradient of that functional. For example, when considering the pure surface minimization problem corresponding to (2.4) without any additional constraint, one postulates that the phasefield evolution is governed by

$$\begin{cases} \frac{\partial \phi}{\partial t} \sim \epsilon \nabla \cdot \frac{\partial a}{\partial \nabla \phi} - \frac{1}{\epsilon} w'(\phi) & \text{in } \Omega, \\ \frac{\partial \phi}{\partial \mathbf{n}} = 0 & \text{on } \partial \Omega, \end{cases} \quad (2.6)$$

i.e. a steepest-descent-type flow. When considering problems where the physically problem is that of maximizing an entropy $\mathcal{S}(\phi)$ instead of minimizing an energy, one would similarly postulate a gradient-ascent-type flow based on \mathcal{S} .

This approach has a number of advantages. Firstly, for a reasonable choice of the proportionality, once any potential steady-state is reached, this state will also be a (local) minimizer of the chosen potential, and will therefore be compatible with the equilibrium formulation. Secondly, it still leaves a large flexibility in the actual description - and therefore the parameterization - of the dynamics, which can be adjusted to fit a large class of problems. Thirdly, the approach is conceptually easy, and, due to the presence of an underlying functional, allows for a relatively simple analysis of the resulting equations. Finally, again due to its close relation with minimization/maximation problems, there is a number of potentially attractive algorithms which, possibly after minor modifications, suggest themselves for the numerical solution at the discrete level.

Despite these advantages, there are also a number of good reasons for deviating from the purely variational setting.

The most obvious is that many physical effects can simply not be properly described based on equilibrium constructs such as a free energy or entropy functional alone. This includes in particular inertial effects, i.e. a resistance to change in time, and therefore a concept which is meaningless for “steady-state” functionals³.

A different (and less obvious) reason within the phasefield approach is that adding “artificial” dynamic effects into the evolution equations can help to better match some experimental or sharp-interface results. A particularly popular example of such a modification is given by **anti-trapping** currents in the context of solidification problems.

After this short summary of the basic idea underlying the phasefield method, the following chapter will outline a number of selected applications in order to clarify some of the typical settings (and the unavoidable challenges) associated with them.

³Even though the dynamics of many of the classical physical equations can be related to other variational principles based on e.g. space-time integral over some action (see e.g. [44]), this need not be the case and in addition need not be well-suited for what one is trying to achieve. Others, such as e.g. the Navier-Stokes equations, do not derive from a classical variational principle.

Chapter 3

An Outline of the Topics Considered

3.1 Pure Phasefield Problems

Even though the ability to capture the effect of surface energies through a smooth description as in Equation (2.4) is theoretically interesting and does already have a number of useful applications, in particular in its multiphase generalization, this by itself is not the primary reason for the increasing popularity of the phasefield method. Instead, it is, as indicated above, the relative ease with which additional effects intervening in combination with surface energy (or entropy) effects can be included into the basic model.

In the simplest cases, this can be achieved by adding e.g. an additional energy contribution $f(\phi, \dots)$ depending only, besides the phasefield values, on a number of **fixed** external parameters. One such example is given by the energy minimization given a fixed undercooling as outlined above in the discussion preceding Equation (2.5). Physically, this corresponds to the case when the solidification is slow enough such that one is able to maintain the material at a constant undercooling, e.g. when one assumes that the material is in contact with an outer thermal reservoir and the heat-conduction within the material happens at a much faster rate than the phase-transformation itself. In this case, an initial solid nucleus will grow, provided that the curvature of its surface is sufficiently low, such that the solidification of the material can liberate enough latent heat to counterbalance the associated increase in the surface energy¹. There are a large number of other physical effects which can be included in much the same manner. These range from in principle relatively simple ones, as for example the influence of gravity for studying the deformations this induces in droplets on a surface, to significantly more complex ones such as the effect of (measured) stored elastic energy contributions on recrystallization processes [78].

Remark 1. Without this ability to include additional “driving forces”, the phasefield method would essentially be reduced to a numerical tool for treating minimal surface problems. This can, even for two phases only, be of some practical interest when considering such problems within geometrically complex regions. There are also a number of more complex questions which can be addressed when one is able to model the interaction of a larger number of phases, as this allows studying e.g. energetically favorable arrangements of bubbles such as in the honeycomb conjecture. Nevertheless, it is clear that this type of problem only accounts for a fraction of the large number of questions being investigated by the phasefield method. \diamond

A big advantage of this type of energy contributions through some interpolation of **fixed** bulk

¹The increase in volume upon an outwards movement, and therefore the rate of energy transformed from a latent to a thermal form during the solidification, is proportional to the surface area $\sim R^2$, whereas the growth in the surface and the associated energy is proportional to the curvature $\sim \frac{1}{R}$. Due to this difference in scalings, the volume increase dominates for a large radius whereas the surface energy dominates when R is small.

energy contributions is that the phasefield is the only actual unknown, and one does therefore not have to consider any of the additional complexities associated with coupled problems, some of which will be outlined in Sections 3.2.1 and 3.2.2. Despite the seeming simplicity, these problems already contain, as far as the phasefield itself is concerned, essentially all the “prototypical ingredients” and challenges also arising in the more complex coupled case.

This will be made use of in Chapter 6 for discussing some of the generic practical issues associated with the numerical treatment of phasefield problems, which are primarily due to the interplay of a and w alone and do not really depend on the precise structure of f . A particular focus there will be on some practical consequences of the use of the obstacle potential for w and the associated constraint of lying within the Gibbs-simplex this imposes on the phasefield values.

This should certainly not be misunderstood in the sense that the choice of f is not important. In fact, one of the major challenges for phasefield models is that large driving forces as compared to the contributions by the basic phasefield term a and w can lead to significant deviations with respect to the expected sharp interface limits, and the majority of Chapter 7 is concerned with more accurate models for the additional energetic contributions to concentration- and/or elasticity-based effects.

This is essentially a question of how one can better model the physics of these additional fields in the presence of diffuse interfaces, and **not** on how the phasefield should react to them. From the point of view of the evolution equation for the phasefield, this is still a problem of precisely the same form, i.e. two contributions due to the bulk and gradient energy (or entropy) densities and some additional driving force - regardless of a potentially increased “internal” complexity - and the primary challenges from a numerical point of view are due to the surface energy contributions and not f .

Remark 2. A different but also very interesting class of problems is obtained when the energy functional (2.4) is not (or not exclusively) supplemented by an additional volumetric contribution but by an additional energy contribution associated with the boundary $\partial\Omega$. In this case, one can consider a part Γ_S of $\partial\Omega$, corresponding to a substrate S , on which wetting phenomena take place. These can be included into the model by adding a contribution of the form ([49] and [9])

$$\int_{\Gamma_S} \Delta\gamma_S h(\phi) \, d\mathbf{s}, \quad (3.1)$$

where $\Delta\gamma_S$ corresponds to the difference of the surface energies between the two phases and the substrate, and h is a suitable interpolation function.

This, together with a gravity term, for example allows studying equilibrium shapes obtained through capillary rise problems within complex geometrical settings but will not be discussed in detail here. \diamond

3.2 Coupled Problems

Assuming a fixed undercooling for the energy functional in Equation (2.5) represents a major simplification. Even though the functional does then depend on the temperature, this temperature is assumed to be given and uniform and is thus essentially just an additional parameter for controlling the energy differences through the amount of undercooling without “actively” participating in the minimization problem. In many cases though, this simplifying assumption cannot be made. In this example, the solidification of the material will lead to a local increase in temperature. If the heat conduction takes place on a time-scale which is comparable with the phase transformation itself or if the material is well-isolated from its surroundings, this cannot (or at least not immediately) be counteracted by giving off heat to the environment, thus invalidating

the assumption of a constant temperature. In order to obtain a physically meaningful model, it therefore becomes necessary to change the role of the temperature to that of an additional unknown.

More generally, a phase transformation process may, besides the temperature, depend on a large number of additional parameters. For more complex solidification problems, it is for example necessary to include additional effects depending on the concentration, which, in most cases, leads to the introduction of an additional unknown. In other processes such as solid-solid transformations, even if performed sufficiently slow to allow for the approximation of a constant temperature, the phase transformation will induce additional stresses within the material, which in turn may significantly influence the transformation process itself (see e.g. [35], [66], [4] and [5] for some recent application-oriented works within the phasefield context). From a modeling point of view, it is again often possible to incorporate these additional effects by complementing the basic phasefield functional with additional mechanical contributions such as the strain energy and an appropriate set of equations for determining the displacement and strains.

3.2.1 Solidification Problems

The Basic Model

As the simplest extension of the situation considered in Section 2, one can consider a setting in which the solidification process is driven solely by an undercooling below the melting temperature of the material under consideration, but where the heat-conduction itself is not fast enough to allow for the approximation of a fixed (external) temperature. In order to capture the dynamics of the temperature field T in a physically meaningful manner, it is necessary to introduce an evolution equation based on the conservation of energy. In the simplest case, this energy can be expressed as a function on ϕ and T alone as $e = e(\phi, T)$, and one postulates

$$\frac{\partial e}{\partial t} = \nabla \cdot \mathbf{q}, \quad (3.2)$$

where the heat-flux \mathbf{q} is usually specified in terms of the material parameters and the temperature.

Alternatively, and in a generalized setting, one may also consider solidification problems which, besides the temperature, also depend on the concentration of one or several chemical components. One such model, with the additional ability to describe an arbitrary number of phases, is proposed in [52] (see also [71]). In this model, the underlying functional is given in terms of the entropy functional

$$\mathcal{S}(\boldsymbol{\phi}, \mathbf{c}, e) = \int_{\Omega} s(\boldsymbol{\phi}, \mathbf{c}, e) - \left(\epsilon a(\boldsymbol{\phi}, \nabla \boldsymbol{\phi}) + \frac{1}{\epsilon} w(\boldsymbol{\phi}) \right) d\mathbf{x} \quad (3.3)$$

where the vectorial functions $\boldsymbol{\phi} : \Omega \rightarrow \mathbf{R}^N$ and $\mathbf{c} : \Omega \rightarrow \mathbf{R}^K$ describe the different phases and components and a and w represent suitable entropy-based generalizations of the surface and bulk energy densities to the multiphase setting. In addition, $\boldsymbol{\phi}$ and \mathbf{c} are locally subject to the constraint $\sum_{\alpha=1}^N \phi^\alpha = 1$ and $\sum_{i=1}^K c_i = 1$, and it is postulated that, based on the conservation of energy and mass and with an appropriate choice of the interaction matrix $\mathbf{L} = (L_{ij})_{0 \leq i \leq K, 0 \leq j \leq K}$,

the energy and concentration follow the evolution equations

$$\begin{aligned}\frac{\partial e}{\partial t} &= -\nabla \cdot \left(L_{00} \nabla \left(\frac{1}{T} \right) + \sum_{j=1}^K L_{0j} \nabla \left(-\frac{\mu_j}{T} \right) \right), \\ \frac{\partial c_i}{\partial t} &= -\nabla \cdot \left(L_{i0} \nabla \left(\frac{1}{T} \right) + \sum_{j=1}^K L_{ij} \nabla \left(-\frac{\mu_j}{T} \right) \right),\end{aligned}$$

where the **chemical potential** $\boldsymbol{\mu}$ is given by $\boldsymbol{\mu} = \frac{\partial f}{\partial \boldsymbol{c}}$. In order to be consistent with the maximization of the entropy functional (3.3), the phasefield is subject to the evolution equation²

$$\tau \epsilon \frac{\partial \phi^\alpha}{\partial t} = \epsilon \left(\nabla \cdot \frac{\partial a}{\partial \nabla \phi^\alpha} - \frac{\partial a}{\partial \phi^\alpha} \right) - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} - \frac{1}{T} \frac{\partial f}{\partial \phi^\alpha} - \Lambda,$$

where Λ is a suitable Lagrange multiplier for the sum-constraint $\sum_{\alpha=1}^N \phi^\alpha = 1$.

Quantitative Phasefield Models

While early phasefield models for purely temperature-driven problems have been quite successful in matching experimental results, it has been observed that the results based on an influence of a concentration vector \boldsymbol{c} - while qualitatively correct - are plagued by a number of artefacts if the interface width is, for numerical reasons, chosen much larger than the physical one ([41], [56] and [19]).

In particular for isothermal solidification problems, a more advanced modeling approach through so-called **quantitative phasefield models** has had, besides a simple increase in computational power, a very significant impact on the ability of phasefield models to more accurately describe the physics of phase transformations. The central idea, initiated by [42], is that the conditions determining the steady-state are based on an equilibration of the temperature and chemical potentials. Whereas models based on simple interpolation procedures of the phase-specific free energy contributions such as $f(\boldsymbol{\phi}, \boldsymbol{c}, T) = \sum_{\alpha=1}^N f^\alpha(\boldsymbol{c}, T) h^\alpha(\boldsymbol{\phi})$ are, with respect to the temperature, therefore expressed in terms of a quantity which is “naturally” expected to be equal between different phases, the same is not true for the additional dependence on the common concentration and leads to excess energy contributions within the interface regions.

A way to significantly reduce these artefacts is, instead of using a single “averaged” concentration field, to introduce an additional internal degree of freedom into the model by “splitting” the given concentration into contributions from each of the coexistent phases, i.e. by introducing **phase-specific concentrations** \boldsymbol{c}^α for each phase α , which, weighted by a suitable interpolation function h^α , satisfy $\sum_{\alpha=1}^N \boldsymbol{c}^\alpha h^\alpha(\boldsymbol{\phi}) = \boldsymbol{c}$. This clearly increases the number of unknowns, and therefore requires introducing additional equations for their specification. As shown in [42], postulating that these phase-specific concentrations \boldsymbol{c}^α are such that the associated chemical potentials corresponding to the f^α are equal allows for a significant improvement of the model. This model has further been analyzed and extended to the multiphase and multicomponent setting by various authors.

In particular, it was realized relatively early by Eiken et al. [25] that this model can also be elegantly interpreted as one where the free energy density is defined as the minimal averaged one $\sum_{\alpha=1}^N f^\alpha(\boldsymbol{c}^\alpha, T) h^\alpha(\boldsymbol{\phi})$ that can be achieved subject to the constraint of the \boldsymbol{c}^α averaging to the total one \boldsymbol{c} . Later works by e.g. [56] and [19] in contrast have primarily focused on a direct description in terms of the chemical potential itself and a very efficient but slightly misleading derivation based on the use of a grand chemical potential approach instead of a free-energy based one. In addition, they replace the evolution equation for \boldsymbol{c} through an in principle equivalent one for $\boldsymbol{\mu}$, which is constructed such that the corresponding evolution of the average concentration

²See the discussion of Equation (3.5) below for an explanation of the term $\frac{1}{T} \frac{\partial f}{\partial \phi^\alpha}$.

as a function of the phasefield and chemical potential matches the standard one through Fick’s law.

The focus in Section 7.1 will be on some primarily practical implications of this type of model. As will be seen there, even though the description in terms of $\boldsymbol{\mu}$ can be a very efficient one for some simple (but practically quite relevant) cases, it is unfortunately much less so in the general case, a point which will again be seen in Section 7.1. In particular, after outlining the model in the spirit of [25] and discussing its relation with the approach in [56] and [19], some practical aspects involved in its actual numerical evaluation - and potential pitfalls - will be discussed. Furthermore, it will be seen that the model generalizes in a straightforward manner to the non-isothermal setting based on a similar definition in terms of an entropy-functional of the form in Equation 3.3, and that these different definitions are consistently linked through the “standard” thermodynamic transforms.

Remark 3. Note that the coupled problems with an increased number of unknown fields and different ways to parameterize the problem give rise to a somewhat delicate question.

As pointed out above, the evolution equation for the phasefield is usually obtained by postulating a gradient-type flow with respect to some functional. This functional is written down in terms of a number of “natural” parameters such as e.g. the phasefield, concentrations, energy or temperature and elastic strains. In order to determine these additional parameters, the evolution equation for ϕ needs to be supplemented with physically appropriate evolution laws for the relevant remaining unknowns (i.e. those which are considered as the independent variables) of the problem. Almost always, there will be interdependencies between the various parameters of the phasefield-functional and the unknowns, either directly by construction (the simpler case), or implicitly. In fact, as noted in e.g. [41] in the solidification context, the steady-state solutions are not truly independent. The question therefore arises how the approach using independent unknowns differs from a **reduced formulation**, where e.g., such as in [41], \mathbf{c} is treated as a function of ϕ (and possibly other parameters) in terms of the steady-state equation it satisfies and whether or not the final result is actually related to a local maximizer or minimizer of the functional.

In the thermodynamic setting, this issue of interrelations between various functions, thermodynamic potentials and parameters is of course well known and in principle also well understood. In particular when adhering to the standard physicist notation of identifying functions and their values, one can easily forget about any possible - then hidden - dependencies. That this notation is nevertheless very useful and does not, at least in many of the more classical situations, cause any serious issues relies fundamentally on an underlying variational structure between the various thermodynamic potentials. This entails that, despite the ubiquitous changes between various interrelated variables, any hidden dependencies have a tendency of simply “dropping out”.

As a simple illustration, consider the example of determining the derivative of an entropy density $s(\phi, \mathbf{c}, e)$ such as in [52] with respect to ϕ . As working directly with the entropy and its (natural) dependence on e can be cumbersome, it is often advantageous to instead rely on a description in terms of the free energy $f(\phi, \mathbf{c}, T)$, in terms of which one has $\frac{\partial s}{\partial \phi} = -\frac{1}{T} \frac{\partial f}{\partial \phi}$.

Formally, this is easily obtained from a partial derivative (i.e. while keeping e constant) of the well-known equality $s = \frac{e-f}{T}$ with respect to ϕ . While this seems logical at first sight, at a closer look this calculation actually raises some questions. Firstly, if f is a function of ϕ , \mathbf{c} and T , then defining s through $\frac{e-f(\phi, \mathbf{c}, T)}{T}$ would make the entropy a function of ϕ , \mathbf{c} , e **and** T . That this is of course not the case is due to the fact that T and e are not independent in this relation. Instead, given (ϕ, \mathbf{c}, e) , there is, under natural assumptions on f , a single $T > 0$ for which the values of $s(\phi, \mathbf{c}, e)$ and $\frac{e-f(\phi, \mathbf{c}, T)}{T}$ coincide. Based on this, one could define T - now as a function of (ϕ, \mathbf{c}, e) - such that this relation always holds, i.e. such that one has

$$s(\phi, \mathbf{c}, e) = \frac{e - f(\phi, \mathbf{c}, T(\phi, \mathbf{c}, e))}{T(\phi, \mathbf{c}, e)}. \quad (3.4)$$

This returns the parameters of s to the original ones, but would seem to invalidate the simple formal differentiation above, as, again under natural assumptions, one a priori actually has

$$\frac{\partial s}{\partial \phi}(\phi, \mathbf{c}, e) = -\frac{1}{T(\phi, \mathbf{c}, e)} \frac{\partial f}{\partial \phi}(\phi, \mathbf{c}, T(\phi, \mathbf{c}, e)) + \frac{\partial}{\partial T} \left(\frac{e - f(\phi, \mathbf{c}, T)}{T} \right) (\phi, \mathbf{c}, T(\phi, \mathbf{c}, e)) \frac{\partial T}{\partial \phi}. \quad (3.5)$$

That the second term actually drops out is due to the fact that s and f are not arbitrary functions but are actually related to one another through a variational property, namely given ϕ , \mathbf{c} and e , s and $T(\phi, \mathbf{c}, e)$ satisfy³

$$s(\phi, \mathbf{c}, e) = \inf_T \left\{ \frac{e - f(\phi, \mathbf{c}, T)}{T} \right\} \quad \text{and} \quad T(\phi, \mathbf{c}, e) = \operatorname{argmin}_\theta \left\{ \frac{e - f(\phi, \mathbf{c}, \theta)}{\theta} \right\},$$

i.e. $T(\phi, \mathbf{c}, e)$ is precisely the T minimizing $\frac{e - f(\phi, \mathbf{c}, T)}{T}$. It is then clear (at least formally) that T is characterized by the Euler-Lagrange equation $\frac{\partial}{\partial T} \frac{e - f(\phi, \mathbf{c}, T)}{T} = 0$, and therefore the relation $\frac{\partial s}{\partial \phi}(\phi, \mathbf{c}, e) = -\frac{1}{T} \frac{\partial f}{\partial \phi}(\phi, \mathbf{c}, T)$ does indeed hold where T is the temperature corresponding to (ϕ, \mathbf{c}, e) .

This has of course been known in a similar form for more than a century (that the derivative is with respect to ϕ instead of e.g. the more classical case of the concentration plays no real role here), but it highlights the fundamental importance that variational interrelations can play for treating potentially interdependent variables as independent⁴.

The same basic principle also provides the answer to the two questions posed above: The formulations in terms of independent variables and in a reduced form are compatible with each other⁵, and the final state will be a local maximizer/minimizer if (and essentially only then!) the steady-state equations allow for an appropriate variational interpretation in terms of the functional. Otherwise, while one can impose independent evolution or steady-state equations on some of the variables, one then either has to explicitly take their dependence on the other parameters into account - a potentially arduous and expensive task - or one loses the relation of the final solution with any constrained (local) maximizer or minimizer of the given functional. \diamond

3.2.2 Solid-Solid Phase Transformation Problems

In recent years, often based on a similar reasoning, a number of quantitative phasefield models have also been proposed for more complex physical models involving - instead of (or in addition to) the concentration fields - the influence of additional elastic energy contributions due to mechanical interactions on phase transformation processes with several solid phases. These additional energy contributions may for example arise from either a prestress of the material, or stresses induced at interfaces due to different crystal structures and/or orientations between different phases and can have a significant influence on the phase-transformation process. Whereas elastic effects within the bulk-phases are relatively well understood, their modeling within an - often artificially large - diffuse interface region in terms of the phasefield approach

³This is similar to the more standard Legendre-Fenchel transform and will be discussed in more detail in Section 5.

⁴In this context, it is important to stress that this is **not** a consequence of (3.4) by itself, which is simply a definition of T for two **given** functions s and f . If s and f are “well-behaved”, one can e.g. invoke an implicit function theorem to obtain information on the derivatives of T based on those of s and f . However, Equation (3.4) does not a priori imply anything particular about $\frac{\partial}{\partial T} \left(\frac{e - f(\phi, \mathbf{c}, T)}{T} \right)$ and could in principle be used for relatively arbitrary functions f and s , provided one can find a function $T(\phi, \mathbf{c}, e)$ such that the equality holds. In particular, it is (obviously) not sufficient to write down an equation similar to (3.4) which “looks like” a Legendre-transformation and then to proclaim independence.

⁵In the sense that, while both approaches might lead to different dynamics and therefore potentially different local minimizers, each will accept the steady-state solution of the other.

has faced similar challenges as the models underlying solidification processes.

Remaining within the small-strain setting, one of the earlier models was based on the use of a single average strain-field ϵ and a “suitable” interpolation of the stiffness-tensor $\mathbf{C}(\phi)$ as e.g. $\mathbf{C}(\phi) = \sum_{\alpha=1}^n \mathbf{C}^\alpha h^\alpha(\phi)$ of the individual stiffness-tensors \mathbf{C}^α with some interpolation function h^α . Using these two quantities, one can naturally define (assuming for simplicity the absence of any prestresses or eigenstrains) a volumetric free energy contribution $f_{el}(\phi, \epsilon) = \frac{1}{2} \epsilon : \mathbf{C}(\phi) : \epsilon$ and the resulting stress tensor $\boldsymbol{\sigma}(\phi, \epsilon) = \mathbf{C}(\phi) : \epsilon = \frac{\partial f_{el}(\phi, \epsilon)}{\partial \epsilon}$.

While this provides a seemingly reasonable interpolation scheme, it was soon realized that its use led to similar artefacts as the ones encountered in the earlier, simpler solidification models. Given the success of the models based on the use of a common chemical potential of all phases instead of the common concentration and based on the analogous roles of $\boldsymbol{\sigma}$ and $\boldsymbol{\mu}$ with respect to the underlying energies, Seinbach and Apel proposed a different scheme in [69] relying on the equality of the stresses instead of that of the strains.

Unfortunately, it has been observed that (see [23]), that, depending on the mechanical setting, both models suffer from excess energy effects within the interface. Durga et al. further observed that these deviations for both models can primarily be attributed to the fact that one does in general neither expect a full equilibration of the strains nor of the stresses. Instead, the corresponding continuity conditions within a sharp interface setting are given by an equality of the normal stresses and an equality of the tangential strains. Based on this, they proposed a model using a “mixed” interpolation approach mimicking this expected behavior within the diffuse interface region, showing a significant improvement in the modeling results. The model was then, within a two-phase setting, further developed and extended in [51], [24] and [64].

In contrast to the simpler mechanical models and the solidification models from Section 3.2.1, these mechanical models, while very elegant in the two-phase case, are unfortunately very difficult to generalize to a multiphase setting. Even though different extensions to the multiphase case have been proposed in e.g. [61], [74], [63] and [62], none of them can be considered fully satisfactory since they either have to rely on a physically undesirable “geometric” simpliciation through a common normal vector between all phase-pairings, or will suffer from a violation of the jump conditions for at least parts of the phases. These difficulties are not unexpected since the multiphase regions in the mechanical model correspond to the intersection of several interfaces in the sharp interface setting - and thus points which are generally associated with singularities in the behavior of the mechanical fields - but it is important to be aware of the limitations and relative advantages and disadvantages associated with each of them. Since the differences in the models are primarily inherited from the particular description chosen in the simpler two-phase case (where, except for the model of [23] and [24], the models coincide) and the chosen description in addition has a significant influence on the computational cost and implementation effort, Section 7.2 will start by summarizing and comparing the various descriptions, both in terms of the formulation itself and in terms of some computational aspects. Subsection 7.2.4 then contains an outline of the different extensions to the multiphase case and some of the issues associated with those. The remainder of Section 7.2 is then devoted to a discussion of some aspects concerning coupled mechanical and chemical calculations, in particular in combination with the jump-condition based mechanical model and the more advanced free energy model from Section 7.1.

Part II

Background

The outline in Chapter 3 clearly indicates the importance of the variational structure and the constraints underlying the phasefield approach in general as well as for the description of the more quantitative models for both the solidification and solid-solid transformations. For this reason, before returning to a more detailed discussion of the applications introduced above, this section will provide some mathematical background on such problems and the closely associated topic of Lagrange multipliers. As this is a topic of both high practical importance and independent mathematical interest, it has been studied in great detail and the literature on the subject is vast. The purpose here is thus not to provide an in-depth discussion, but just to outline some basics which will be used in the sequel. This will be complemented by numerous literature references where a more detailed presentation and related but more advanced results may be found.

The central difference between constrained minimization problems and their unconstrained counterparts is that both the characterization of and the search for (local) minima is made more difficult by the fact that, in the former case, it is not sufficient to focus just on some differential information of the objective function alone. Instead, it is obvious that the constraints must also be a part of the formulation of any necessary and/or sufficient condition describing minima and of any algorithm aiming at their determination. As a first step, the following Chapter 4 contains a discussion of equality and inequality-constrained problems. In both cases, the focus is on the description of first-order necessary conditions for local minimizers, which naturally leads to the notion of Lagrange-multipliers.

The results obtained for both the equality- and inequality constrained settings can be conveniently summarized using the concept of Lagrange-functionals, which will be discussed in Section 5.1. These on the one hand allow for a simple formal derivation of first-order necessary and second-order necessary resp. sufficient conditions for constrained local minimizers. On the other hand, they are also the basis for both the design and the analysis of a number of algorithms for the numerical treatment of such problems. Section 5.2 then gives some background on an a priori different topic, namely the Legendre-Fenchel transform, which has a fundamental role in thermodynamics.

Chapter 4

Equality and Inequality Constrained Problems

4.1 Linear and Nonlinear Equality Constraints - the Finite-Dimensional Case

Linear Equality Constraints The simplest setting for constrained minimization problems is that of minimizing a smooth (e.g. C^2) function f defined on \mathbb{R}^n subject to $m \leq n$ linear equality constraints

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{b}_i \cdot \mathbf{x} = c_i, \quad i = 1, \dots, m. \end{cases} \quad (4.1)$$

The question to be answered is then under what conditions one can assert that a point $\bar{\mathbf{x}}$ is a (local) minimizer for the problem (4.1).

The first obvious point is that $\bar{\mathbf{x}}$ needs to satisfy $\mathbf{b}_i \cdot \bar{\mathbf{x}} = \mathbf{b}_i^T \bar{\mathbf{x}} = c_i$ for $i = 1, \dots, m$ (and thus in particular the constraints need to be consistent). In order for such a point $\bar{\mathbf{x}}$ to be local minimizer, the value $f(\bar{\mathbf{x}})$ has to be no larger than those of $f(\mathbf{x})$ with \mathbf{x} in a neighborhood of $\bar{\mathbf{x}}$ which, in addition, also satisfy $\mathbf{b}_i \cdot \mathbf{x} = c_i, i = 1, \dots, m$. Due to this restriction and the linearity of the constraints, any admissible variation $\delta\mathbf{x}$ of $\bar{\mathbf{x}}$ has to satisfy $\mathbf{b}_i \cdot \delta\mathbf{x} = 0, i = 1, \dots, m$. Combining the equality-constraints into a vector-equation by inserting the vectors \mathbf{b}_i row-wise into the matrix

$$\mathbf{B} = \begin{pmatrix} \mathbf{b}_1^T \\ \mathbf{b}_2^T \\ \dots \\ \mathbf{b}_m^T \end{pmatrix},$$

this can be written in a more compact form as $\mathbf{B} \cdot \delta\mathbf{x} = \mathbf{0}$ or $\delta\mathbf{x} \in \text{Ker}(\mathbf{B})$. As all admissible points are of the form $\bar{\mathbf{x}} + \delta\mathbf{x}$ with $\delta\mathbf{x} \in \text{Ker}(\mathbf{B})$, the optimality condition is thus given by $f(\bar{\mathbf{x}} + \delta\mathbf{x}) \geq f(\bar{\mathbf{x}})$ for all sufficiently small $\delta\mathbf{x} \in \text{Ker}(\mathbf{B})$.

It remains to convert this into a differential characterization. Given any $\delta\mathbf{x} \in \text{Ker}(\mathbf{B})$, the vector $\mathbf{x} = \bar{\mathbf{x}} + t\delta\mathbf{x}$ clearly satisfies $\mathbf{B}\mathbf{x} = \mathbf{c}$ provided $\bar{\mathbf{x}}$ does so, i.e. is admissible, and, since $\bar{\mathbf{x}}$ is assumed to be a local minimizer, will satisfy $f(\mathbf{x}) \geq f(\bar{\mathbf{x}})$ or

$$0 \leq \frac{f(\mathbf{x}) - f(\bar{\mathbf{x}})}{t} = \frac{f(\bar{\mathbf{x}} + t\delta\mathbf{x}) - f(\bar{\mathbf{x}})}{t}$$

for sufficiently small t . Letting $t \rightarrow 0$ and since f is assumed smooth¹, it then follows that any minimizer $\bar{\mathbf{x}}$ has to satisfy the **first-order necessary condition (FONC)**

$$\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \geq 0 \quad \forall \delta \mathbf{x} \in \text{Ker}(\mathbf{B}). \quad (4.2)$$

As $\text{Ker}(\mathbf{B})$ is a linear subspace of \mathbb{R}^n , this condition can be further simplified to the equality²

$$\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} = 0 \quad \forall \delta \mathbf{x} \in \text{Ker}(\mathbf{B}). \quad (4.3)$$

In contrast to the unconstrained case where $\delta \mathbf{x}$ can be chosen arbitrarily, this does not imply $\nabla f(\bar{\mathbf{x}}) = \mathbf{0}$. Instead, condition (4.3) only states that $\nabla f(\bar{\mathbf{x}})$ lies in the orthogonal complement of the kernel of \mathbf{B} ,

$$\nabla f(\bar{\mathbf{x}}) \in \text{Ker}(\mathbf{B})^\perp.$$

In order to obtain a more explicit condition on $\nabla f(\bar{\mathbf{x}})$, it is thus sufficient to give a characterization of $\text{Ker}(\mathbf{B})^\perp$. Each of the equations $\mathbf{b}_i \cdot \delta \mathbf{x} = 0$ characterizing the kernel of \mathbf{B} is nothing but the description of a hyperplane with the vector \mathbf{b}_i acting as the normal vector. In the case of a single constraint, it is geometrically obvious that $\text{Ker}(\mathbf{b}_i^T)^\perp = \text{Span}(\mathbf{b}_i)$ (see Figure 4.1), or, equivalently, if \mathbf{b}_i is considered as the linear mapping $\mathbf{b}_i : \mathbb{R} \ni \lambda \rightarrow \lambda \mathbf{b}_i \in \mathbb{R}^n$, that $\text{Ker}(\mathbf{b}_i^T)^\perp = \text{Range}(\mathbf{b}_i)$.

More generally, the kernel of \mathbf{B} is just the intersection of such hyperplanes, $\text{Ker}(\mathbf{B}) = \cap_i \text{Ker}(\mathbf{b}_i^T)$. As $\cap_{i=1}^m \text{Ker}(\mathbf{b}_i^T) \subset \text{Ker}(\mathbf{b}_i^T)$ and $A \subset B$ implies $A^\perp \supset B^\perp$ for any subsets A, B of \mathbb{R}^n , it is clear that $\text{Ker}(\mathbf{B})^\perp \supset \text{Ker}(\mathbf{b}_i^T)^\perp$, $i = 1, \dots, m$, and thus also³

$$\text{Ker}(\mathbf{B})^\perp \supset \text{Span}(\{\text{Ker}(\mathbf{b}_i^T)^\perp\}_{1 \leq i \leq m}) = \text{Span}(\{\mathbf{b}_i\}_{1 \leq i \leq m}) = \text{Range}(\mathbf{B}^T).$$

A restatement of the converse conclusion $\text{Ker}(\mathbf{B})^\perp \subset \text{Range}(\mathbf{B}^T)$ is simply that any vector which is not in $\text{Range}(\mathbf{B}^T)$ is not orthogonal to $\text{Ker}(\mathbf{B})$. In fact, any vector $\mathbf{x} \notin \text{Range}(\mathbf{B}^T)$ can be written in the form $\mathbf{x} = \mathbf{y} + \sum_{i=1}^m \alpha_i \mathbf{b}_i$ where⁴ $\mathbf{y} \neq \mathbf{0}$ is orthogonal to $\text{Range}(\mathbf{B}^T)$ (apply e.g. a Gram-Schmidt orthogonalization procedure using \mathbf{x} and the $\{\mathbf{b}_i\}_{1 \leq i \leq m}$). This means precisely that $\mathbf{y} \perp \mathbf{b}_i$, $i = 1, \dots, m$ or $\mathbf{y} \in \text{Ker}(\mathbf{B})$, and thus $\mathbf{y} \cdot \mathbf{x} = \|\mathbf{y}\|^2 \neq 0$ shows that $\mathbf{x} \notin \text{Ker}(\mathbf{B})^\perp$.

In combination with the same argument applied to $\text{Ker}(\mathbf{B}^T)$ and $\text{Range}(\mathbf{B})$, this implies the following fundamental theorem:

Theorem 1. (*Fundamental theorem of linear algebra*)

Let $\mathbf{B} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $m, n \geq 1$ be any real matrix. Then the orthogonal decompositions

$$\mathbb{R}^n = \text{Ker}(\mathbf{B}) \oplus \text{Range}(\mathbf{B}^T) \quad (4.4)$$

and

$$\mathbb{R}^m = \text{Ker}(\mathbf{B}^T) \oplus \text{Range}(\mathbf{B})$$

hold.

Using this result, $\nabla f(\bar{\mathbf{x}}) \in \text{Ker}(\mathbf{B})^\perp$ is thus equivalent to $\nabla f(\bar{\mathbf{x}}) \in \text{Range}(\mathbf{B}^T)$, i.e. there exists a vector $\boldsymbol{\lambda} \in \mathbb{R}^m$ of **Lagrange-multipliers** such that⁵

$$-\nabla f(\bar{\mathbf{x}}) = \mathbf{B}^T \boldsymbol{\lambda}. \quad (4.5)$$

¹Clearly, Gâteaux-differentiability is sufficient here.

²In fact, otherwise assuming that $\delta \mathbf{x}$ is a variation such that $\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} > 0$, $-\delta \mathbf{x}$ is necessarily also a legitimate variation, but, with $\nabla f(\bar{\mathbf{x}}) \cdot (-\delta \mathbf{x}) = -\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} < 0$, contradicting the optimality of $\bar{\mathbf{x}}$.

³This is again intuitively clear in \mathbb{R}^n , as this simply expresses that any linear combination of the vectors $\{\mathbf{b}_i\}_{1 \leq i \leq m}$ is orthogonal to all vectors which are orthogonal to all of the $\{\mathbf{b}_i\}_{1 \leq i \leq m}$.

⁴As will be seen below, the analogue of this seemingly simple statement in \mathbb{R}^n in terms of a separation theorem is the major source of difficulty for an analogous result in the infinite-dimensional setting. The basic difficulty in the latter case is that, while the range of any continuous linear operator is always a subspace, this subspace (contrary to the situation in \mathbb{R}^n) need not be closed. An element x not in the range of the operator might thus lie in its closure and can therefore be approximated arbitrarily close by elements in the range.

⁵Note that using $-\nabla f$ instead of ∇f is simply a matter of convenience here as one could just as well change the sign of $\boldsymbol{\lambda}$.

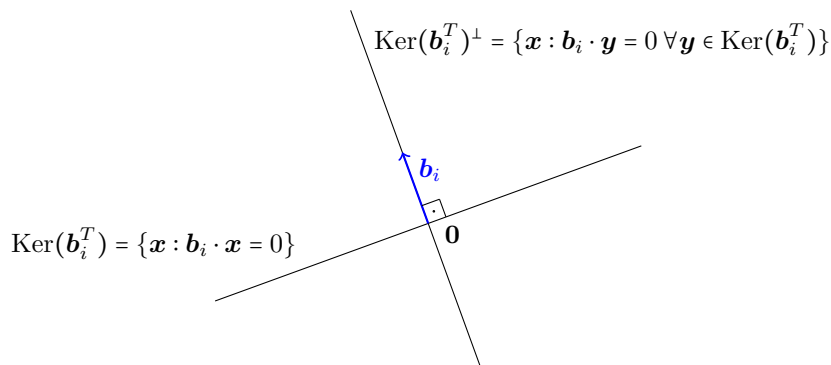


Figure 4.1: Relationship between $\text{Ker}(\mathbf{b}_i^T)$ and $\text{Ker}(\mathbf{b}_i^T)^\perp$

Remark 4. Despite its simplicity, the result above is, as the name suggests, fundamental for finite-dimensional linear problems (constrained or not) as it gives a very convenient decomposition of the domain and image spaces of any linear operator.

Among the many consequences, one should mention the well-known **Fredholm alternative**: The equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ is solvable iff \mathbf{b} is orthogonal to all solutions \mathbf{y} of the homogeneous transposed equation $\mathbf{A}^T\mathbf{y} = \mathbf{0}$, which is just a reformulation of the second statement $\text{Range}(\mathbf{A}) = \text{Ker}(\mathbf{A}^T)^\perp$ above.

In particular, the operator $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is surjective iff $\text{Ker}(\mathbf{A}^T) = \{\mathbf{0}\}$, i.e. \mathbf{A}^T is injective. As all (if any) solutions of an equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ only differ by elements lying in $\text{Ker}(\mathbf{A})$, a similar conclusion is that a linear operator $\mathbf{A} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is injective (or $\text{Ker}(\mathbf{A}) = \{\mathbf{0}\}$) iff \mathbf{A}^T is surjective (i.e. $\text{Range}(\mathbf{A}^T) = \mathbb{R}^m$) as the first statement then shows that $\mathbb{R}^m = \text{Ker}(\mathbf{A}) \oplus \mathbb{R}^m$. \diamond

Remark 5. The preceding arguments are prototypical for the more difficult situations considered below.

First of all, under mild assumptions, they can be applied essentially without change in the nonlinear case considered in the next paragraph after a simple linearization.

In addition, these result can in many cases of practical interest be extended to the infinite-dimensional setting where the matrix \mathbf{B} above is replaced by a suitable linear operator acting between two Banach spaces $\mathbf{X} \rightarrow \mathbf{Y}$. The only fundamental change here (at least for smooth functionals $\mathcal{F} : X \rightarrow \mathbb{R}$) is that this requires some additional topological assumptions, under which the decisive Theorem 1 can be generalized in the form of the **closed range theorem** (see Thm. 2 below).

Finally, the treatment of inequality constrained problems of the form $\mathbf{B}\mathbf{x} \leq \mathbf{c}$ for some given matrix \mathbf{B} and a vector \mathbf{c} also runs much along the same lines. A major difference in this case is that, while the Euler-Lagrange equation characterizing any local minimizers could be simplified to an equality in the setting above, this is no longer possible when inequalities are present. This leads to the requirement of replacing the subspace defining the kernel of \mathbf{B} through an intersection of the hyperplanes by a cone defined through a suitable intersection of half-spaces. \diamond

Nonlinear Equality Constraints While most of the equality constraints arising in the applications considered in this thesis are indeed linear, this is of course not always the case. Fortunately, under relatively mild additional hypothesis on the equality constraints, the FONC for the more general case where \mathbf{x} is subject to $m \leq n$ nonlinear equality constraints

$$h_i(\mathbf{x}) = c_i, \quad i = 1, \dots, m \quad (4.6)$$

can be deduced from the previous consideration after a simple linearization of the h_i , $i = 1, \dots, m$. More precisely, one would like to replace the admissible variations $\delta\mathbf{x} \in \text{Ker}(\mathbf{B})$ at a feasible

point $\bar{\mathbf{x}}$ from the FONC of the previous section with the vectors $\delta\mathbf{x} \in \text{Ker}(\mathbf{h}'(\bar{\mathbf{x}})')$, where \mathbf{h} is the (column) vector composed of the $h_i(\mathbf{x})$, $\mathbf{h}(\mathbf{x}) = \begin{pmatrix} h_1(\mathbf{x}) & h_2(\mathbf{x}) & \dots & h_m(\mathbf{x}) \end{pmatrix}^T$. Provided this is legitimate, one can apply the theory from the linear case to deduce the existence of a (not necessarily unique) Lagrange multiplier $\boldsymbol{\lambda} \in \mathbb{R}^m$ such that any local minimizer satisfies

$$\nabla f(\bar{\mathbf{x}}) = (\mathbf{h}'(\bar{\mathbf{x}}))^T \boldsymbol{\lambda} = \nabla \mathbf{h}(\bar{\mathbf{x}}) \boldsymbol{\lambda}.$$

Before entering into a more detailed discussion, it is instructive to consider a simple counterexample in \mathbb{R}^2 which shows that this is not always possible. Consider a smooth function $f(\mathbf{x})$ with \mathbf{x} constrained to lie in the admissible set consisting of all points on the vertical axis, but with this (in its most natural form) linear constraint artificially made nonlinear by rewriting this set using the constraint $h(\mathbf{x}) = x^2 = 0$. It is clear from the discussion in the previous section that at any minimizer, there will be a $\lambda \in \mathbb{R}$ such that $\nabla f = \lambda \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, i.e. the only non-vanishing component of the gradient is along the x -direction.

Linearizing h though, any admissible point with $x = 0$ will satisfy

$$\nabla h(\mathbf{x}) = h'(\mathbf{x})^T = \begin{pmatrix} x \\ 0 \end{pmatrix}^T = \mathbf{0}.$$

It is obvious that this vector is useless for expressing the FONC above.

Remark 6. Roughly speaking, one can see that the difficulty in the above example arises as there are two different “linearizations” involved, which in this case do not lead to the same conclusions. One linearization is the more geometric one based on the admissible set \mathcal{A} (in this case the y -axis) itself and the translated tangent space $T_{\mathcal{A}}(\mathbf{x})$ to this set obtained by considering the tangents to all smooth curves lying in \mathcal{A} and passing through \mathbf{x} (in this case again the y -axis). The other one is a more algebraic one based on the linearization of the function \mathbf{h} describing the admissible set and the set of directions lying in $\text{Ker}(\mathbf{h}'(\mathbf{x}))$. As seen in the example above, this set corresponds to the y -axis for the choice $h(\mathbf{x}) = x$ and to the whole of \mathbb{R}^2 (i.e. a significantly larger set) for the choice $h(\mathbf{x}) = x^2$. In contrast to $T_{\mathcal{A}}(\mathbf{x})$, this set thus depends on the particular expression chosen for \mathbf{h} and may or may not coincide with the former.

While it is intuitively to be expected that one is interested in the first set for the constrained problems considered here, it is also clear that the second one is much easier to deal with.

The principal question to be answered is therefore under which conditions both subsets coincide. These so-called **constraint qualification conditions** appear generally in constrained optimization problems where the definition of the admissible set \mathcal{A} involves systems of equalities and inequalities. The difficulty is always due to the above dichotomy, where one has, on the one hand, a purely geometrically defined tangent set to \mathcal{A} which appears naturally in the Euler-Lagrange equation characterizing the minimizer, and, on the other hand, an appropriate linearization of the constraint equations which one would like to combine with Lagrange multipliers in order to obtain a more explicit representation of this set. \diamond

As for the linearly constrained case, the FONCs at any purported local minimizer $\bar{\mathbf{x}}$ are again based on considering nearby points also satisfying the inequality constraints. For all such points sufficiently close to $\bar{\mathbf{x}}$, one has, by the assumption on $\bar{\mathbf{x}}$, that $f(\mathbf{x}) - f(\bar{\mathbf{x}}) \geq 0$. In order to obtain a first-order necessary condition based on directional derivatives, one would now like to divide by t and take the limit of $t \rightarrow 0$ for points of the form $\bar{\mathbf{x}} + t\delta\mathbf{x}$ with an admissible direction $\delta\mathbf{x}$. There is a slight difficulty with this approach, as, due to the potential curvature of the admissible set, it is in general not possible to approach $\bar{\mathbf{x}}$ along a straight line while remaining within the admissible set \mathcal{A} , i.e. it is not exactly legitimate to compare the values of $f(\mathbf{x})$ along a line of the type above with those of $\bar{\mathbf{x}}$. Instead, one has to compare the value of $\bar{\mathbf{x}}$ with points lying “almost” on such a straight line, i.e. points $\mathbf{x} \in \mathcal{A}$ which can be written in the form $\mathbf{x} = \bar{\mathbf{x}} + t\delta\mathbf{x} + o(t)$ for some direction $\delta\mathbf{x}$. In fact, provided f is sufficiently smooth and there is some sequence t_n

tending to 0 such that an associated sequence $\mathbf{x}_n = \bar{\mathbf{x}} + t_n \delta \mathbf{x} + o(t_n)$ lying in \mathcal{A} exists, one then has

$$0 \leq \frac{f(\mathbf{x}_n) - f(\bar{\mathbf{x}})}{t_n} = \frac{f(\bar{\mathbf{x}}) + t \nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} + o(t_n) - f(\bar{\mathbf{x}})}{t_n} = \nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} + o(1)$$

and thus, in the limit, $\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \geq 0$.

Denoting the ‘‘convergence in \mathcal{A} ’’ with $\mathbf{x}_n \xrightarrow{\in \mathcal{A}} \bar{\mathbf{x}}$, i.e. all points \mathbf{x}_n (and $\bar{\mathbf{x}}$) lie in \mathcal{A} , the set of directions for which this construction is possible is given by the tangential set

$$T_{\mathcal{A}}(\bar{\mathbf{x}}) = \left\{ \delta \mathbf{x} \in \mathbb{R}^n : \exists t_n \rightarrow 0, \mathbf{x}_n \xrightarrow{\in \mathcal{A}} \bar{\mathbf{x}} \text{ such that } \delta \mathbf{x} = \lim_{t_n \rightarrow 0} \frac{\mathbf{x}_n - \bar{\mathbf{x}}}{t_n} \right\} \quad (4.7)$$

called the **sequential/Bouligand cone**. With this set, the FONC in this more general case is given by

$$\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \geq 0 \quad \forall \delta \mathbf{x} \in T_{\mathcal{A}}(\bar{\mathbf{x}}). \quad (4.8)$$

According to the preceding discussion, in order to obtain a description in terms of Lagrange multipliers as outlined above, one now needs to ensure that $T_{\mathcal{A}}(\bar{\mathbf{x}})$ is in fact a subspace (allowing to simplify the inequality into an equality) and that this subspace is given by the kernel of $\mathbf{h}'(\bar{\mathbf{x}})$ in order to conclude that $\nabla f(\bar{\mathbf{x}}) \in \text{Ker}(\mathbf{h}')^\perp = \text{Range}((\mathbf{h}'(\bar{\mathbf{x}}))^T)$ and thus that

$$-\nabla f(\bar{\mathbf{x}}) = (\mathbf{h}'(\bar{\mathbf{x}}))^T \boldsymbol{\lambda} = \nabla \mathbf{h}(\bar{\mathbf{x}}) \boldsymbol{\lambda} \quad (4.9)$$

for some $\boldsymbol{\lambda} \in \mathbb{R}^m$.

A first observation in this sense (also recall the example above) is that, if \mathbf{h} is smooth, e.g. $\mathbf{h} \in C^1(\mathbb{R}^n)$, then $T_{\mathcal{A}}(\bar{\mathbf{x}}) \subset \text{Ker}(\mathbf{h}'(\bar{\mathbf{x}}))$ as, for any direction $\delta \mathbf{x} \in T_{\mathcal{A}}(\bar{\mathbf{x}})$ and the associated sequences \mathbf{x}_n and t_n , one has

$$\mathbf{0} = \frac{\mathbf{h}(\mathbf{x}_n) - \mathbf{c}}{t_n} = \frac{\mathbf{h}(\bar{\mathbf{x}} + t_n \delta \mathbf{x} + o(t_n)) - \mathbf{c}}{t_n} = \frac{\mathbf{h}(\bar{\mathbf{x}}) + t_n \nabla \mathbf{h}'(\bar{\mathbf{x}}) \delta \mathbf{x} + o(t_n) - \mathbf{c}}{t_n} = \mathbf{h}'(\bar{\mathbf{x}}) \delta \mathbf{x} + o(1)$$

and thus in the limit $\delta \mathbf{x} \in \text{Ker}(\mathbf{h}'(\bar{\mathbf{x}}))$.

All that can a priori be said about vectors $\delta \mathbf{x} \in \text{Ker}(\mathbf{h}'(\bar{\mathbf{x}}))$ is that

$$\mathbf{h}(\bar{\mathbf{x}} + t \delta \mathbf{x}) = \mathbf{h}(\bar{\mathbf{x}}) + t \mathbf{h}'(\bar{\mathbf{x}}) \delta \mathbf{x} + o(t) = \mathbf{c} + o(t), \quad (4.10)$$

i.e. they satisfy the equality constraint up to an error of $o(t)$. What is required in order to show that $\delta \mathbf{x}$ lies in $T_{\mathcal{A}}$ though is not that the equation defining \mathcal{A} is satisfied up to an order of $o(t)$, but that there is a point \mathbf{x} within a distance of $o(t)$ to $\bar{\mathbf{x}} + t \delta \mathbf{x}$ actually satisfying the equation. This conclusion based on Equation (4.10) requires a stability result ensuring that the $o(t)$ -error in the defining equation for \mathcal{A} can be compensated by an equal-order correction to the point $\bar{\mathbf{x}} + t \delta \mathbf{x}$. In other words, one has to ensure that \mathbf{h} is a local homeomorphism at $\bar{\mathbf{x}}$, i.e. a continuous mapping with a continuous inverse. This is a classical question though, for which it is known (see e.g. [27]) that this is the case iff $\mathbf{h}'(\bar{\mathbf{x}} + t_n \delta \mathbf{x})$ is surjective, which is in turn guaranteed by $\mathbf{h}'(\bar{\mathbf{x}})$ being surjective.

4.2 Linear and Nonlinear Equality Constraints - the Infinite-Dimensional Case

Even though the setting considered in the previous section already contains the essential ideas underlying the Lagrange multiplier approach for equality-constrained problems, it is not strictly speaking applicable to the applications outlined above. Instead, one has to consider a more general situation in which the domain of the function(al) to be minimized is given by an appropriate subset of a function space and thus generally infinite-dimensional. In addition, while it may happen - as in the case of the minimization with a prescribed volume - that there is only a finite number of equality constraints, there are other situations in which this need not be the case anymore⁶.

Both aspects entail some additional difficulties absent from the considerations in the previous section. Even though the discretization of any of the examples considered in this section ultimately reduces the situation to the previous finite-dimensional one, the underlying continuous description should not be completely forgotten.

Remark 7. This is not just a matter of theoretical interest but also has implications from a “practical” point of view.

On the one hand, despite the ultimately discrete nature of the computational problems, a continuous description almost always leads to significantly shorter calculations as compared to an analogous one at the discrete level, and is in particular also heavily made use of when deriving e.g. the phasefield equation and other related ones from the underlying functional. Even though the engineering community enjoys some additional liberties in this respect as compared to the mathematical one (which will also be taken here), it can still be helpful to have an idea of which calculations are potentially problematic and which are not. In particular, it is clear that the continuous description will become increasingly important as one increases the resolution in the discrete case, and an ill-defined description at the continuous level is then likely to cause issues at the discrete one.

On the other hand, a basic understanding of some important relations at the continuous level - such as e.g. the gradient and negative divergence being related in a “transpose-like” fashion - allows for a relatively simple intuitive interpretation of some a priori quite complex questions at the discrete level. \diamond

The remainder of this section will therefore provide a quick sketch of how and to what degree the arguments from Section 4.1 can be carried over to the function space setting. The central result replacing the use of the fundamental theorem of linear algebra in the previous section is a very similar characterization of the decomposition of the domain and range space through the closed range theorem 2 below, with the role of the matrix \mathbf{B} and its transpose \mathbf{B}^T in Theorem 1 being replaced by the operator \mathcal{B} defining the constraints and its adjoint \mathcal{B}^* . Provided the conditions of the theorem are satisfied, one can use this characterization to derive, based on essentially identical arguments as in the finite-dimensional case, the existence of an appropriate Lagrange multiplier for the constraint and the way it enters into the first-order necessary condition through \mathcal{B}^* .

Remark 8. The closed range theorem is a fundamental theorem of functional analysis and can therefore be found in essentially any introductory text on this topic. The outline below mostly borrows from [16], [48],[80] and [20], to which the reader is referred for more details. More in-depth descriptions and generalizations can be found in e.g. [14], [37] and [57]. \diamond

⁶A particularly important example is given by the incompressibility condition $\nabla \cdot \mathbf{u} = 0$ encountered when dealing with fluid flow problems.

Linear equality constraints Consider first the case when $\mathcal{F} : X \rightarrow \mathbb{R}$ is a (smooth) functional defined on some Banach or Hilbert space⁷ X and the equality constraints are specified through a (possibly unbounded⁸) linear operator $\mathcal{B} : \text{Dom}(\mathcal{B}) \subset X \rightarrow Y$,

$$\mathcal{B}x = y,$$

where Y is a second Banach space.

Remark 9. The simplest case is the one where Y is finite-dimensional, corresponding to a finite number of (linear) equality constraints, specified through m continuous linear functionals $b_i^* \in X^*$,

$$\langle b_i^*, x \rangle_{X^*, X} = c_i, \quad i = 1, \dots, m.$$

For the phasefield, one may for example have $X = H^1(\Omega) \cap L^\infty(\Omega)$, $\mathcal{F} = \mathcal{E}_\epsilon$ and $\int_\Omega \phi \, dx = V$. For the (Navier-)Stokes equation one instead has to consider the more general case where X may e.g. be a subspace of $H^1(\Omega)$ and Y equal to $L^2(\Omega)$. \diamond

As in the previous section, standard differential calculus implies that the derivative $\mathcal{F}'(\bar{x})$ of \mathcal{F} (as an element of X^*) at any minimizer \bar{x} must satisfy

$$\langle \mathcal{F}'(\bar{x}), \delta x \rangle_{X^*, X} \geq 0 \quad \forall \delta x \in \text{Ker}(\mathcal{B}),$$

and then again, as $\text{Ker}(\mathcal{B})$ is a linear subspace of X ,

$$\langle \mathcal{F}'(\bar{x}), \delta x \rangle_{X^*, X} = 0 \quad \forall \delta x \in \text{Ker}(\mathcal{B}). \quad (4.11)$$

Defining⁹ the **annihilator** M^a of a subset of M as the set of all $x^* \in X^*$ such that $\langle x^*, x \rangle_{X^*, X} = 0$ for all $x \in M$, one thus has $\mathcal{F}'(\bar{x}) \in \text{Ker}(\mathcal{B})^a$. As in the case of \mathbb{R}^n , one would now like to, whenever this is possible, provide a more explicit representation of $\text{Ker}(\mathcal{B})^a$. This requires introducing, in analogy to the transposed matrix, the **adjoint operator** $\mathcal{B}^* : \text{Dom}(\mathcal{B}^*) \subset Y^* \rightarrow X^*$:

Definition 1. (Adjoint operator, see e.g. [16], p. 43f)

Let $\mathcal{B} : \text{Dom}(\mathcal{B}) \subset X \rightarrow Y$ be a linear operator which is densely defined¹⁰. The domain of \mathcal{B}^* is defined to be the set of $y^* \in Y^*$ such that there exists a constant c such that

$$|\langle y^*, \mathcal{B}(x) \rangle_{Y^*, Y}| \leq c \|x\|_X \quad \forall x \in \text{Dom}(\mathcal{B}),$$

i.e. $\text{Dom}(\mathcal{B}^*)$ is the (linear) subspace of Y^* for which the linear mapping $D(\mathcal{B}) \ni x \mapsto g(x) = \langle y^*, \mathcal{B}x \rangle$ is uniformly bounded. As $\text{Dom}(\mathcal{B})$ is by assumption dense in X , g can be extended by continuity to a unique bounded linear operator on all of X , allowing to identify g with an element of X^* . This association in turn defines a linear mapping associating with each $y^* \in \text{Dom}(\mathcal{B}^*)$ an element of X^* , which will be denoted by $\mathcal{B}^* y^*$ and by definition satisfies

$$\langle y, \mathcal{B}x \rangle_{Y^*, Y} = \langle \mathcal{B}^* y^*, x \rangle_{X^*, X} \quad \forall y^* \in \text{Dom}(\mathcal{B}^*), x \in \text{Dom}(\mathcal{B}). \quad (4.12)$$

⁷I.e. a complete normed vector space, resp. one where the norm can additionally be derived from an inner product.

⁸Meaning that \mathcal{B} need not be continuous on the whole space. Typical examples for such operators are differentiation operators between various spaces, which need not be bounded on the whole space but are so on a dense subset of sufficiently smooth functions.

⁹The reason for introducing a slightly “modified version” of the orthogonal complement here is that in the more general Banach space setting, not every linear subspace admits an orthogonal complement. When X is a Hilbert space, this is the case though for every **closed** subspace. As the kernel of every linear continuous operator is always closed, in this case one can thus always replace $\text{Ker}(\mathcal{B})^a$ with $\text{Ker}(\mathcal{B})^\perp$ since both notions coincide then (see e.g. [16]).

¹⁰Meaning its definition is such that it “makes sense” on a dense subset of functions in X .

The adjoint operator is thus constructed such as to extend the defining property of the transposed matrix $\mathbf{B}^T : \mathbb{R}^m \rightarrow \mathbb{R}^n$ with respect to the Euclidian scalar product in \mathbb{R}^m resp. \mathbb{R}^n ,

$$\mathbf{y} \cdot (\mathbf{B}\mathbf{x}) = (\mathbf{B}^T \mathbf{y}) \cdot \mathbf{x},$$

to the more general setting given by the duality product between the spaces above. The hope is of course to obtain an analogue of the fundamental theorem of linear algebra for this case. Unfortunately, this is not possible without some additional topological assumptions, which hold in many cases of practical interest though.

Their statement requires another definition (see e.g. [16], p. 43):

Definition 2. (Graphs and closed maps)

The **graph** $\text{graph}(\mathcal{B})$ of a map (linear or not) $\mathcal{B} : X \rightarrow Y$ is the subset of $X \times Y$ given by

$$\text{graph}(\mathcal{B}) = \{(x, y) \in X \times Y : y = \mathcal{B}(x)\}.$$

\mathcal{B} is **closed** if $\text{graph}(\mathcal{B})$ is a closed subset of $X \times Y$ (equipped with the standard product topology)¹¹.

With this definition, the fundamental result for the characterization of $\text{Ker}(\mathcal{B})^a$ is contained in the following (see e.g. thm. 2.16 [16], thm. 2.13 [48], §7.5 [79])

Theorem 2. (Closed range theorem)

Let $\mathcal{B} : \text{Dom}(\mathcal{B}) \subset X \rightarrow Y$ be an unbounded linear operator which is densely defined and closed. Then the following are equivalent:

1. $\text{Range}(\mathcal{B})$ is closed.
2. $\text{Range}(\mathcal{B}^*)$ is closed.
3. $\text{Range}(\mathcal{B}) = \text{Ker}(\mathcal{B}^*)^a$.
4. $\text{Range}(\mathcal{B}^*) = \text{Ker}(\mathcal{B})^a$.

Remark 10. Note that there are two different notions of “closedness” involved in this theorem. The first one - that of being a closed map - depends on the simultaneous convergence of a sequence of points (x_n) in X and its image sequence $(y_n) = \mathcal{B}x_n$ in Y and requires that, if both converge, the limit in Y is the image of the limit in X . The second one - that of having a closed range - is more concerned with the image space, i.e. given an arbitrary sequence (y_n) in Y which converges in Y , the limit has to lie in $\text{Im}(\mathcal{B})$, meaning there only has to be some point in X such that $y = \mathcal{B}x$ (but which, a priori, need not be related to any particular sequence in X). \diamond

Assuming that the conditions of this theorem are satisfied for \mathcal{B} , one can finally conclude from (4.11) that at a local minimizer \bar{x} ,

$$-\mathcal{F}'(\bar{x}) \in \text{Ker}(\mathcal{B})^a = \text{Range}(\mathcal{B}^*),$$

or, more explicitly, that there exists a Lagrange multiplier $\lambda \in Y^*$ such that

$$-\mathcal{F}'(\bar{x}) = \mathcal{B}^* \lambda. \tag{4.13}$$

Remark 11. An important example for the use of this theorem is e.g. the role of the pressure gradient in relation with the divergence-free constraint in fluid dynamics. If one considers a sufficiently smooth domain Ω and the (negative) divergence operator $-\text{div}$ as an operator from $\mathbf{H}_0^1(\Omega) = (H_0^1(\Omega))^n$ into $L^2(\Omega)$, then it is actually a continuous operator between these

¹¹Here this reduces to: For any sequence $x_n \in D(\mathcal{B})$ such that $x_n \rightarrow x \in X$ and $y_n := \mathcal{B}x_n \rightarrow y \in Y$, it must hold that $x \in D(\mathcal{B})$ and $y = \mathcal{B}x$

spaces (since $|\sum_{i=1}^N \frac{\partial u}{\partial x_i}|^2 \leq n \sum_{i=1}^N |\frac{\partial u}{\partial x_i}|^2$ and thus $\|-\operatorname{div}(\mathbf{u})\|_{L^2(\Omega)} \leq \sqrt{n}\|\mathbf{u}\|_{H_0^1(\Omega)}$) and therefore defined on all of $\mathbf{H}_0^1(\Omega)$. By Schwarz's inequality, one further has $\|-\operatorname{div}(u)q\|_{L^2(\Omega)} \leq \|-\operatorname{div}(u)\|_{L^2(\Omega)}\|q\|_{L^2(\Omega)}$ and thus the mapping $\mathbf{u} \mapsto \int_{\Omega} -\operatorname{div}(\mathbf{u})q \, d\mathbf{x}$ is uniformly bounded for all q in $L^2(\Omega)$ with constant $c = \sqrt{n}\|q\|_{L^2(\Omega)}$. The domain of its adjoint $-\operatorname{grad} : L^2(\Omega) \rightarrow \mathbf{H}^{-1}(\Omega) = \mathbf{H}^{-1}(\Omega)$, where $\mathbf{H}^{-1}(\Omega) = (H^{-1}(\Omega))^n$ is the dual of $\mathbf{H}_0^1(\Omega)$, will therefore be all of $L^2(\Omega)$. For smooth functions, Green's formula shows that

$$\int_{\Omega} -\operatorname{div}(\mathbf{u})q \, d\mathbf{x} = \int_{\Omega} \mathbf{u} \cdot \nabla q \, d\mathbf{x} - \int_{\partial\Omega} q\mathbf{u} \cdot \mathbf{n} \, ds = \int_{\Omega} \mathbf{u} \cdot \nabla q \, d\mathbf{x}$$

since \mathbf{u} vanishes on the boundary and the adjoint is therefore an extension of the standard gradient operator¹². It further follows that the ‘‘Green-type’’ formula

$$\langle \operatorname{grad} q, \mathbf{u} \rangle_{\mathbf{H}^{-1}(\Omega), H_0^1(\Omega)} = \int_{\Omega} (-\operatorname{div} \mathbf{u})q \, d\mathbf{x}$$

holds, regardless of whether the duality pairing on the left may be written as an integral and whether an integration by parts may be justified or not.

The assumption for Theorem 2 can be shown to hold (see e.g. [73]) and thus $\operatorname{Ker}(\operatorname{div})^a = \operatorname{Range}(-\operatorname{grad}) = \operatorname{Range}(\operatorname{grad})$. The crucial implication in this case is that any functional $\mathbf{f} \in \mathbf{H}^{-1}(\Omega)$ which vanishes on all divergence-free functions in $H_0^1(\Omega)$ can be written as the gradient of a scalar function $p \in H_0^1(\Omega)$, i.e. $\langle \mathbf{f}, \mathbf{v} \rangle_{\mathbf{H}^{-1}, H_0^1(\Omega)} = 0$ for all $\mathbf{v} \in \mathbf{H}_0^1(\Omega)$ implies that $\mathbf{f} = \nabla p$, $p \in L^2(\Omega)$. ◊

Remark 12. The definition of the adjoint and Theorem 2 are also, similar to Remark 4, the central ‘‘ingredients’’ for e.g. the very useful Fredholm alternative (see e.g. theorem 2.27 in [48]). ◊

Nonlinear Equality Constraints As in the finite-dimensional case, the theory developed for the case of linear equality constraints carries over to nonlinear equality constraints provided some regularity condition holds. This is summarized in the following result:

Theorem 3. (*Tangents and normals to a level set (Graves-Lyusternik theorem), thm. 5.35 [20]*)

Let X and Y be Banach spaces and let \mathcal{A} be given by

$$\mathcal{A} = \{u \in X : h(x) = 0\},$$

where the map $h : X \rightarrow Y$ is continuously differentiable near $x \in \mathcal{A}$. If $h'(x)$ is surjective, then $T_{\mathcal{A}}(x)$ and $N_{\mathcal{A}}(x)$ are the linear subspaces given by

$$T_{\mathcal{A}}(x) = \{\delta x \in X : \langle h'(x), \delta x \rangle = 0\}, \quad N_{\mathcal{A}}(x) = (h'(x))^* Y^*$$

and $T_{\mathcal{A}}(x) = N_{\mathcal{A}}(x)^{\circ}$, $N_{\mathcal{A}}(x) = T_{\mathcal{A}}(x)^{\circ}$.

This is essentially the same conclusion as in the finite-dimensional case, i.e. if the derivative of the nonlinear mapping is surjective (recall the counterexample in Remark 6), the only admissible variations are the ones on which the linearization of the constraint vanishes, and the derivative of the function(al) vanishing on all these directions then implies that there is some multiplier such that

$$\mathcal{F}'(x) + (h'(x))^* \lambda = 0.$$

¹²More precisely, even though the expression $\operatorname{grad} q$ above does not a priori make ‘‘classical’’ sense for $q \in L^2(\Omega)$, the image of q under the gradient operator is, by the construction of the adjoint, simply **defined** as the unique element in $\mathbf{g} \in \mathbf{H}^{-1}(\Omega)$ satisfying $\langle \mathbf{g}, \mathbf{u} \rangle_{\mathbf{H}^{-1}(\Omega), H_0^1(\Omega)} = \int_{\Omega} q(-\operatorname{div}(\mathbf{u})) \, d\mathbf{x} \forall \mathbf{u} \in \mathbf{H}_0^1(\Omega)$.

Remark 13. Note that the surjectivity assumption in Theorem 3 implicitly subsumes what is required for Theorem 2. In fact, by the assumption that h is C^1 , h' is defined everywhere on X (and not just densely). Furthermore, since $h'(x)$ is assumed surjective, the graph of $h'(x)$ is simply the product space $X \times Y$, and therefore obviously a closed subset of $X \times Y$ (its complement being the empty set, which is always open). Similarly, $\text{Range}(h'(x))$ equals the whole space Y and is therefore again closed, from which the remaining conclusions follow immediately. \diamond

4.3 Linear and Nonlinear Inequality Constraints

In many applications, one not only has to deal with equality constraints but with possible additional inequality constraints. The primary example in the phasefield context is of course the box-constraints $0 \leq \phi \leq 1$ that may be imposed on the phasefield, or, in combination with an equality constraint in the multi-phase case, the restriction to the Gibbs-simplex, $0 \leq \phi^\alpha \leq 1$, $\sum_{\alpha=1}^N \phi^\alpha = 1$. Another physically important example is given by elasto-plastic applications, in which the admissible stresses are assumed to be restricted by a yield criterion of the form $f(\boldsymbol{\sigma}) \leq \sigma^y$, with σ^y possibly depending on additional internal parameters.

In contrast to the equality constrained case, which leads to first-order optimality conditions in terms of orthogonality to certain linear subspaces as in Equations (4.3) and (4.11), the Euler-Lagrange equations for inequality constrained problems lead to first-order conditions in the form of inequalities. This in turn leads to a characterization involving cones instead of subspaces. For this reason, it is necessary to find a suitable replacement for the fundamental theorem of linear algebra (resp. the closed graph theorem), here primarily in the form of Farka's lemma.

As for the equality constraints, the first case considered will be the finite-dimensional one in combination with linear inequality constraints, a setting which already contains the essential ideas. It will then shortly be sketched how these can be extended to some more complex situations.

Linear Inequality Constraints in the Finite-Dimensional Case Consider again the minimization of a smooth objective function f defined on \mathbf{R}^n , but now subjected to $m \leq n$ linear inequality constraints:

$$\begin{cases} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{b}_i \cdot \mathbf{x} \leq c_i, \quad i = 1, \dots, m. \end{cases} \quad (4.14)$$

Just as for the equality constrained case, the FONC for a local minimizer $\bar{\mathbf{x}}$ satisfying the constraints is given by the Euler-Lagrange equation

$$\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \geq 0 \quad (4.15)$$

with $\delta \mathbf{x}$ ranging over the ‘‘admissible’’ directions. The fundamental difference with respect to the previous situation lies precisely in the form of this set, which is now not simply given by the kernel of the matrix \mathbf{B} composed of the \mathbf{b}_i . Instead, it is clear that there is actually more freedom in choosing the $\delta \mathbf{x}$ as $\text{Ker}(\mathbf{B})$ is now only a subset of the admissible variations¹³. In fact, for each inequality constraint one has to distinguish two possible situations:

1. $\bar{\mathbf{x}}$ satisfies $\mathbf{b}_i \cdot \bar{\mathbf{x}} < c_i$, i.e. the i -th constraint is **inactive**. By continuity, the same will hold for any \mathbf{x} sufficiently close to $\bar{\mathbf{x}}$ or equivalently for any $\bar{\mathbf{x}} + \delta \mathbf{x}$ with $|\delta \mathbf{x}|$ sufficiently small. In this case, at least locally, the constraint effectively does not impose any actual restriction on the problem and can be disregarded.
2. $\bar{\mathbf{x}}$ satisfies $\mathbf{b}_i \cdot \bar{\mathbf{x}} = c_i$, i.e. the i -th constraint is **active** or **saturated**. In this case, the only admissible variations are such that $\mathbf{b}_i \cdot \delta \mathbf{x} \leq 0$, which, instead of a hyperplane, now specifies an entire half-space.

¹³If $\bar{\mathbf{x}}$ satisfies $\mathbf{B}\bar{\mathbf{x}} \leq \mathbf{c}$, $\bar{\mathbf{x}} + \delta \mathbf{x}$ obviously also satisfies $\mathbf{B}(\bar{\mathbf{x}} + \delta \mathbf{x}) = \mathbf{B}\bar{\mathbf{x}} + \mathbf{B}\delta \mathbf{x} = \mathbf{B}\bar{\mathbf{x}} \leq \mathbf{c}$ for any $\delta \mathbf{x} \in \text{Ker}(\mathbf{B})$.

Partitioning the constraints into the inactive ones,

$$\mathcal{I}_I(\bar{\mathbf{x}}) = \{i \in 1, \dots, n : \mathbf{b}_i \cdot \bar{\mathbf{x}} < c_i\} \quad (4.16)$$

and the active ones,

$$\mathcal{I}_A(\bar{\mathbf{x}}) = \{i \in 1, \dots, n : \mathbf{b}_i \cdot \bar{\mathbf{x}} = c_i\}, \quad (4.17)$$

the set of admissible variations is therefore given by the set $T_K(\bar{\mathbf{x}})$ of $\delta\mathbf{x}$ such that

$$T_K(\bar{\mathbf{x}}) = \{\delta\mathbf{x} \in \mathbb{R}^n : \mathbf{b}_i \cdot \delta\mathbf{x} \leq 0 \quad \forall i \in \mathcal{I}_A(\bar{\mathbf{x}})\},$$

i.e. an intersection of (closed) half-spaces.

Definition 3. (Cones and polar cones)

- A **cone** (also sometimes called a **pointed cone**) is any subset C such that, if $\mathbf{x} \in C$, then $\alpha\mathbf{x} \in C$ for all $\alpha \geq 0$.
- Given any subset M of \mathbb{R}^n , the **(negative) polar cone** M^- is the set of vectors \mathbf{y} such that

$$\mathbf{y} \cdot \mathbf{x} \leq 0 \quad \forall \mathbf{x} \in M.$$

Similarly, its negative, the **dual cone** (or **positive polar cone**) is the set $M^+ = -M^-$ defined by

$$\mathbf{y} \cdot \mathbf{x} \geq 0 \quad \forall \mathbf{x} \in M.$$

With this definition, it is clear¹⁴ that the admissible variations $T_K(\bar{\mathbf{x}})$ form a cone.

The optimality condition (4.15) corresponds, by the very definitions, to $f'(\bar{\mathbf{x}}) \in T_K^+$ or, using the more common notation, to

$$-\nabla f(\bar{\mathbf{x}}) \in T_K^-(\bar{\mathbf{x}}), \quad (4.18)$$

which is thus the conclusion replacing $\nabla f(\bar{\mathbf{x}}) \in \text{Ker}(\mathbf{B})^\perp$ from the equality-constrained case. In order to obtain a more explicit formulation, it remains to replace the characterization $\text{Ker}(\mathbf{B})^\perp = \text{Range}(\mathbf{B}^T)$ by an analogous statement describing $T_K^-(\bar{\mathbf{x}})$. This is the purpose of the following

Lemma 1. (*Farka's lemma, cor. 2.29 [58]*)

Let \mathbf{B} be an $m \times n$ -matrix and let

$$K = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{B}\mathbf{x} \leq \mathbf{0}\}.$$

Then (see also Figure 4.2a)

$$K^- = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y} = \mathbf{B}^T \boldsymbol{\mu}, \boldsymbol{\mu} \in \mathbb{R}^m, \boldsymbol{\mu} \geq \mathbf{0}\}. \quad (4.19)$$

Proof. As for the fundamental theorem of linear algebra, one implications is a simple consequence of the definition of the transposed matrix¹⁵. In fact, denoting the set on the right-hand side of (4.19) by \tilde{K} , one has

$$(\mathbf{x}, \mathbf{B}^T \boldsymbol{\mu})_{\mathbb{R}^m} = (\mathbf{B}\mathbf{x}, \boldsymbol{\mu})_{\mathbb{R}^n} \leq 0 \quad \forall \mathbf{x} \in K, \boldsymbol{\mu} \in \tilde{K},$$

as $\mathbf{B}\mathbf{x}$ has only non-positive and $\boldsymbol{\mu}$ only non-negative entries by the definitions of K and \tilde{K} , i.e. $\tilde{K} \subset K^-$. For the other implication, observe that \tilde{K} is a closed convex cone. If there is a $\mathbf{z} \in K^-$

¹⁴Either by “geometrical insight” or by noting that the defining condition $\mathbf{b}_i \cdot \delta\mathbf{x} \leq 0$ is stable under multiplication by non-negative scalars.

¹⁵ There it is the conclusion that $\text{Range}(\mathbf{B}^T)\boldsymbol{\mu} \subset \text{Ker}(\mathbf{B})^\perp$, as for any $\mathbf{x} \in \text{Ker}(\mathbf{B})$ $\mathbf{x} \cdot \mathbf{B}^T \boldsymbol{\mu} = \mathbf{B}\mathbf{x} \cdot \boldsymbol{\mu} = \mathbf{0}$. In contrast, the analogous statement to the other conclusion $\text{Ker}(\mathbf{B})^\perp \subset \text{Range}(\mathbf{B}^T)$ is trickier as it requires a separation property similar to the one used here.

such that $z \notin \tilde{K}$, z can therefore be strictly separated from \tilde{K} , i.e. there is some vector $l \in \mathbb{R}^n$ such that

$$l \cdot z > \alpha \geq l \cdot y \quad \forall y \in \tilde{K}.$$

As each $y \in \tilde{K}$ is of the form $B^T \mu$ with $\mu \geq 0$, one thus has

$$Bl \cdot \mu \leq \alpha \quad \forall \mu \geq 0,$$

which is only possible if $Bl \leq 0$ (otherwise just let the entry in μ corresponding to a positive entry in Bl tend to $+\infty$). In addition, as $\mu = 0$ is admissible, this implies $\alpha \geq 0$. Since $Bl \leq 0$ is just the same as saying $l \in K$, one would therefore have $z \cdot l > \alpha \geq 0$, but this contradicts $z \in K^-$. \square

Remark 14. While the proof of Lemma 1 may seem somewhat technical through the use of a separation theorem, the geometric idea is actually rather simple. The separation theorem and the first part of the theorem imply that there is a half-space containing $\tilde{K} \subset K^-$ and, in addition, the existence of a vector in the **other** half-space, which, just as all elements of \tilde{K} , forms an angle of at least 90° with any vector in K . As \tilde{K} always contains the vectors b_i , i.e. the ones which are precisely orthogonal to one of the subspaces defining K , there is no way of arranging a half-plane such that it contains both \tilde{K} and that a vector on the other side does not make an acute angle with K (see Figure 4.2b). \diamond

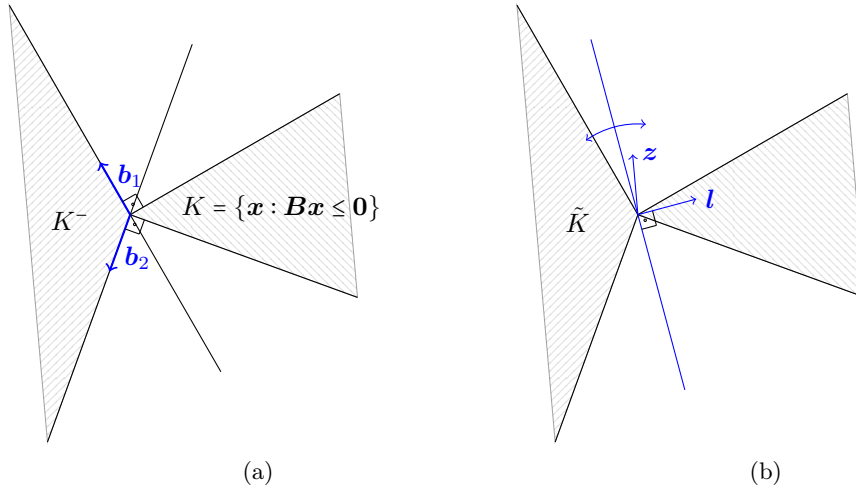


Figure 4.2: The (negative) polar cone in Farkas lemma (a) and the geometrical idea underlying its proof (b).

Applying the lemma to $-\nabla f(\bar{x})$ with the matrix $B_{\mathcal{I}_A(\bar{x})}$ constructed (row-wise) from the vectors b_i^T with $i \in \mathcal{I}_A(\bar{x})$ replacing B leads to

$$-\nabla f(\bar{x}) = B_{\mathcal{I}_A(\bar{x})}^T \mu_{\mathcal{I}_A(\bar{x})} = \sum_{i \in \mathcal{I}_A(\bar{x})} \mu_i b_i,$$

where each of the μ_i appearing above has to be non-negative, $\mu_i \geq 0, i \in \mathcal{I}_A(\bar{x})$.

As the notation involving the index set $\mathcal{I}_A(\bar{x})$ is somewhat cumbersome and only meaningful in the discrete case, it is common to use an equivalent formulation which is obtained by noting that $B_{\mathcal{I}_A(\bar{x})}^T \mu_{\mathcal{I}_A(\bar{x})} = B^T \mu$, with B the analogue of $B_{\mathcal{I}_A(\bar{x})}$ but containing all the vectors b_i , provided all $\mu_i, i \in \mathcal{I}_I(\bar{x})$ are set to zero. These are precisely the indices though for which $c_i - b_i \cdot \bar{x} \geq 0$, whereas the indices in $\mathcal{I}_A(\bar{x})$, i.e. the set where the μ_i may be non-zero, are given by

$c_i - \mathbf{b}_i \cdot \bar{\mathbf{x}} = 0$. Both conditions can compactly be combined by introducing the **complementarity conditions** $\mu_i(c_i - \mathbf{b}_i \cdot \bar{\mathbf{x}}) = 0$, $i = 1, \dots, m$. Finally, as both μ_i and $c_i - \mathbf{b}_i \cdot \bar{\mathbf{x}}$ are required to be non-negative and thus the sum $\sum_{i=1}^m \mu_i(c_i - \mathbf{b}_i \cdot \bar{\mathbf{x}})$ can only vanish if all terms are zero, the characterization of $\nabla f(\bar{\mathbf{x}})$ can equivalently be expressed by

$$\begin{cases} -\nabla f(\bar{\mathbf{x}}) = \mathbf{B}^T \boldsymbol{\mu}, \\ \mathbf{c} - \mathbf{B}\bar{\mathbf{x}} \geq \mathbf{0}, \\ \boldsymbol{\mu} \geq \mathbf{0}, \\ \boldsymbol{\mu} \cdot (\mathbf{c} - \mathbf{B}\bar{\mathbf{x}}) = 0. \end{cases} \quad (4.20)$$

Nonlinear Inequality Constraints in the Finite-Dimensional Case As for the equality-constrained case, the conclusions for the linear inequalities can be carried over to the nonlinear case through a simple linearization of the (active) constraints provided some constraint qualification conditions hold. More precisely, for a constraint-set $\mathcal{A} = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, i = 1, \dots, m\}$ prescribed through $m \leq n$ inequality constraint, the first-order necessary condition is, as in Equation (4.8), that $\nabla f(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \geq 0$ for all $\delta \mathbf{x} \in T_{\mathcal{A}}$ i.e. $-\nabla f(\bar{\mathbf{x}}) \in T_{\mathcal{A}}(\bar{\mathbf{x}})^-$ where the Bouligand cone $T_{\mathcal{A}}(\bar{\mathbf{x}})$ is defined as in Equation (4.7). What one would like to conclude from this is that any minimizer $\bar{\mathbf{x}}$ subject to $g(\mathbf{x}) \leq \mathbf{0}$ is of the form

$$-\nabla f(\bar{\mathbf{x}}) = \sum_{i \in \mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})} \mu_i \nabla g_i(\bar{\mathbf{x}}) \quad (4.21)$$

with $\mu_i \geq 0$, where $\mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})$ denotes the set of active constraints, i.e. the indices for which $g_i(\bar{\mathbf{x}}) = 0$. All this requires is to show that the Bouligand cone appearing in the characterization of $\nabla f(\bar{\mathbf{x}})$ for any local minimizer is the same as the linearizing cone

$$L(\bar{\mathbf{x}}, K) := \{\delta \mathbf{x} : \nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \leq 0, i \in \mathcal{I}_{\mathcal{A}}\}$$

consisting of those directions which, in first order, do not lead to an increase in the values of the g_i which are at their maximal admissible value 0 at $\bar{\mathbf{x}}$. If such is the case, the conclusion follows immediately from Farka's lemma, as one then has $T_{\mathcal{A}}(\bar{\mathbf{x}}) = \{\delta \mathbf{x} : \mathbf{B}\delta \mathbf{x} \leq \mathbf{0}\}$ with $\mathbf{B} = (\nabla \mathbf{g}_{\mathcal{I}_{\mathcal{A}}}(\bar{\mathbf{x}}))^T$ and by Lemma 1 $T_{\mathcal{A}}(\bar{\mathbf{x}})^- = \{\mathbf{y} : \mathbf{y} = \mathbf{B}^T \boldsymbol{\mu}, \boldsymbol{\mu} \geq \mathbf{0}\} = \{\mathbf{y} : \mathbf{y} = \nabla \mathbf{g}_{\mathcal{I}_{\mathcal{A}}}(\bar{\mathbf{x}}) \boldsymbol{\mu}, \boldsymbol{\mu} \geq \mathbf{0}\}$ and thus together with the first order necessary condition in Equation (4.21).

Linking the two cones is again a question of being able to ensure that, any $\delta \mathbf{x}$ in $L(\bar{\mathbf{x}}, K)$ can actually be obtained as the limit of the directions showing towards points actually lying in \mathcal{A} . This will hold if for any $\delta \mathbf{x}$ such that $\nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} \leq 0$ for all $i \in \mathcal{I}_{\mathcal{A}}$, one can find an actually admissible point \mathbf{x} in \mathcal{A} such that $\delta \mathbf{x} = \mathbf{x} - \bar{\mathbf{x}} + o(\|\delta \mathbf{x}\|)$ and is therefore again a question of stability.

There are a variety of different conditions ensuring this (see e.g. [58], [13] or [27] for a more detailed discussion). A relatively strong condition ensuring this is the linear independence constraint qualification condition, requiring that the gradients $\nabla g_i(\bar{\mathbf{x}})$ for $i \in \mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})$ be linearly independent. A weaker condition also based on the gradients themselves is the Mangasarian-Fromovitz constraint qualification condition, requiring that there be some vector \mathbf{d} such that $\nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} < 0$ for all $i \in \mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})$. Another popular condition, which in addition applies even when the g_i are not necessarily differentiable, is Slater's constraint qualification condition, which requires for the g_i , $i \in \mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})$ to be convex functions together with the existence of some point \mathbf{y} such that $\mathbf{g}_{\mathcal{I}_{\mathcal{A}}(\bar{\mathbf{x}})}(\mathbf{y}) < \mathbf{0}$.

Remark 15. Note that, with respect to any active constraint i , this is not an issue for a direction $\delta \mathbf{x}$ such that $\nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} < 0$ as the differentiability of g_i automatically implies that $g_i(\bar{\mathbf{x}} + t\delta \mathbf{x}) < 0$ for t sufficiently small, i.e. the constraint will definitely be satisfied as $\bar{\mathbf{x}}$ is approached along this direction. The only potential difficulties therefore arises for those $\delta \mathbf{x}$ for which $\nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} = 0$. As for such directions, one always has (regardless of any qualification condition) $g_i(\bar{\mathbf{x}} + t\delta \mathbf{x}) =$

$g_i(\bar{\mathbf{x}}) + t\nabla g_i(\bar{\mathbf{x}}) \cdot \delta \mathbf{x} + o(t) = o(t)$, one can ensure feasibility of the a very nearby (i.e. $o(t)$ close) point by a “tiny nudge” in any direction such that $\nabla g_i(\bar{\mathbf{x}}) \cdot \mathbf{d} < 0$, provided this does not lead to a violation of the remaining active inequality constraints. The basic idea underlying the conditions above is to ensure the existence of a direction \mathbf{d} such that this is always possible. \diamond

Remark 16. The results above can be generalized quite significantly in various directions. Firstly, one can replace the condition $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ with more abstract conditions. Secondly, many of the ideas also carry over to the infinite-dimensional setting under some additional topological assumptions, with the adjoint again replacing the transpose matrix. This will not be discussed in detail here, and the reader is referred to e.g. [37] and [14] for a discussion in relation with Lagrange multipliers or e.g. [22] for some discussion on generalization of Farka’s lemma above. A very well-written introduction to some closely related background can also be found in [7], in particular concerning various extensions of classical results from functional analysis to situations involving convex cones instead of the entire spaces. \diamond

Chapter 5

Lagrange Functions and Legendre-Fenchel Duality

5.1 Lagrange Functions

The results for the minimization of $f(\mathbf{x})$ subject to the equality-constraint $\mathbf{h}(\mathbf{x}) = \mathbf{c}$ and $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ from Section 4.1 can conveniently be combined and summarized by introducing the **Lagrange function** or **Lagrangian**

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda} \cdot (\mathbf{h}(\mathbf{x}) - \mathbf{c}) + \boldsymbol{\mu} \cdot \mathbf{g}(\mathbf{x}). \quad (5.1)$$

In fact, it can be observed that, provided one restrict $\boldsymbol{\mu}$ to be non-negative, simply imposing that the gradient $\nabla_{\mathbf{x}}$ of L with respect to \mathbf{x} should vanish,

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \nabla f(\mathbf{x}) + \nabla \mathbf{h}(\mathbf{x}) \boldsymbol{\lambda} + \nabla \mathbf{g}(\mathbf{x}) \boldsymbol{\mu} \stackrel{!}{=} \mathbf{0}, \quad (5.2)$$

one recovers the same form of relation between the gradient of f and the multipliers which was obtained in the previous chapter. In addition, again demanding that the gradient of L with respect to $\boldsymbol{\lambda}$ should vanish,

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{h}(\mathbf{x}) - \mathbf{c} \stackrel{!}{=} \mathbf{0}, \quad (5.3)$$

one recovers the equality constraint. Finally, imposing that the $\nabla_{\boldsymbol{\mu}} L$ should be non-positive leads to

$$\nabla_{\boldsymbol{\mu}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}) \stackrel{!}{\leq} \mathbf{0} \quad (5.4)$$

and therefore the original inequality constraint on $\mathbf{g}(\mathbf{x})$.

Even though this has, up to this point, no deeper meaning, it is obvious that provided that this is a very convenient “mnemonic” for obtaining the correct structure of the first-order optimality conditions derived to a relatively lengthy argument based on orthogonality relations resp. relations between cones and their polar cones in the previous section. As will be sketched below, there is in fact a close relation between Lagrangian and the original constrained optimization problem and the usefulness of the Lagrangian goes far beyond providing a simple means of “guessing” the correct optimality conditions.

The basis for this approach is contained in the following observations: Whenever the equality constraint is satisfied for some \mathbf{x} , the term in $\boldsymbol{\lambda}$ vanishes. Since L does then not depend on $\boldsymbol{\lambda}$ for such an \mathbf{x} and taking the supremum over all $\boldsymbol{\lambda}$ thus has no effect, $\sup_{\boldsymbol{\lambda} \in \mathbb{R}^m} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu} \cdot \mathbf{g}(\mathbf{x})$. Secondly, a similar argument can be done in terms of $\boldsymbol{\mu}$ and $\mathbf{g}(\mathbf{x})$. If \mathbf{x} is such that $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$, the term $\boldsymbol{\mu} \cdot \mathbf{g}(\mathbf{x})$ is necessarily less than or equal to zero. If some entry $g_i(\mathbf{x})$ is strictly negative, the supremum over all $\mu_i \geq 0$ is obviously achieved for $\mu_i = 0$. In

constrast, if $g_i(\mathbf{x}) = 0$, L does not depend on μ_i and any μ_i will therefore do for achieving the supremum $\sup_{\mu_i \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$. Both cases together therefore show that if \mathbf{x} satisfies the inequality constraint, taking the supremum over the $\boldsymbol{\mu} \geq \mathbf{0}$ selects $\boldsymbol{\mu}$ such that the second part $\boldsymbol{\mu} \cdot \mathbf{g}(\mathbf{x}) = 0$ of the complementarity conditions holds and, if this is the case, also eliminates the last term of the Lagrangain in Equation (5.1).

Whenever \mathbf{x} satisfies both constraints, one therefore has

$$\sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}). \quad (5.5)$$

In contrast, when the equality constraint is not satisfied and thus $\mathbf{h}(\mathbf{x}) - \mathbf{c} \neq \mathbf{0}$ there exists some $\boldsymbol{\lambda} \in \mathbb{R}^m$ such that $\boldsymbol{\lambda} \cdot (\mathbf{h}(\mathbf{x}) - \mathbf{c}) > 0$, and thus by “scaling” $\boldsymbol{\lambda}^*$, one has $\sup_{\boldsymbol{\lambda} \in \mathbb{R}^m} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = +\infty$. Similarly, if any entry of $\mathbf{g}(\mathbf{x})$ is greater than 0, it suffices to let the corresponding component μ_i tend to $+\infty$ to show that $\sup_{\boldsymbol{\mu} \geq \mathbf{0}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ is then also $+\infty$.

In summary, taking the supremum in $\boldsymbol{\lambda}$ and $\boldsymbol{\mu} \geq \mathbf{0}$ to define the **primal function** $L_P(\mathbf{x})$, this function satisfies

$$L_P(\mathbf{x}) := \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \begin{cases} f(\mathbf{x}) & \mathbf{h}(\mathbf{x}) - \mathbf{c} = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \\ +\infty & \text{else,} \end{cases} \quad (5.6)$$

and, since the infinite values are obviously of no interest with respect to a minimization, it follows further that the original problem of minimizing $f(\mathbf{x})$ subject to the constraints can equivalently be expressed in terms of the **primal problem**

$$\inf_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x}) = \inf_{\mathbf{x} \in \mathbb{R}^n} \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}). \quad (5.7)$$

While this minimization problem is now in principle a free minimization problem in \mathbf{x} , this is an essentially formal difference, which, in this form, is hard to put to any practical use. The idea is instead to look at an - a priori different - problem obtained by exchanging the order of the inf and sup, i.e. by instead considering the **dual problem**

$$\sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} \inf_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) \quad (5.8)$$

where the **dual function** is defined by

$$L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) := \inf_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) + \boldsymbol{\lambda} \cdot (\mathbf{h}(\mathbf{x}) - \mathbf{c}) + \boldsymbol{\mu} \cdot \mathbf{g}(\mathbf{x})\}. \quad (5.9)$$

Remark 17. A first point to be noted is that the **dual problem**, while still a problem potentially subject to the constraint $\boldsymbol{\mu} \geq \mathbf{0}$ has, at least in terms of $L_D(\boldsymbol{\lambda}, \boldsymbol{\mu})$, a significantly simpler structure than the original minimization problem for $f(\mathbf{x})$. In fact, whereas \mathbf{x} is potentially subject two both a nonlinear equality constraint and a nonlinear inequality constraint, there is no constraint on $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ is solely subject to the simple constraint $\boldsymbol{\mu} \geq \mathbf{0}$.

The latter observation is of course only really advantageous if the inequality constraints on the original variable \mathbf{x} are of a more complex form. This is for example not the case for the restriction of the phasefield to the Gibbs-simplex which will be discussed in more detail in Chapter 6, since the inequality restrictions on the primal unknown are also of the form $\boldsymbol{\phi} \geq \mathbf{0}$. Nevertheless, the practical difficulty of the restriction to the Gibbs-simplex is not primarily due to the inequality constraints themselves (this could be handled using a simple truncation), but their coupling through an additional sum-constraint, an aspect which never occurs in a dual problem. \diamond

A further indication why the dual problem may be useful is given by the following two simple observations, which do not require any assumptions on $f(\mathbf{x})$, $\mathbf{h}(\mathbf{x})$ or $\mathbf{g}(\mathbf{x})$:

Lemma 2. (see e.g. chapter 4 [58])

The dual function is concave and one always has

$$\sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \inf_{\boldsymbol{x} \in \mathbb{R}^n} L_P(\boldsymbol{x}). \quad (5.10)$$

Proof. Firstly, $f(\boldsymbol{x}) + \boldsymbol{\lambda} \cdot (\mathbf{h}(\boldsymbol{x}) - \mathbf{c}) + \boldsymbol{\mu} \cdot \mathbf{g}(\boldsymbol{x})$ is linear and thus also (even though “border-line”) concave in both $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$. As the infimum of an arbitrary family of concave functions is concave, the first claim follows.

By the definition of the supremum, it is also clear that

$$L_P(\boldsymbol{x}) = \sup_{\{(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\mu}}); \tilde{\boldsymbol{\mu}} \geq \mathbf{0}\}} \{f(\boldsymbol{x}) + \tilde{\boldsymbol{\lambda}} \cdot (\mathbf{h}(\boldsymbol{x}) - \mathbf{c}) + \tilde{\boldsymbol{\mu}} \cdot \mathbf{g}(\boldsymbol{x})\} \geq f(\boldsymbol{x}) + \boldsymbol{\lambda} \cdot (\mathbf{h}(\boldsymbol{x}) - \mathbf{c}) + \boldsymbol{\mu} \cdot \mathbf{g}(\boldsymbol{x})$$

regardless of the choice of $\boldsymbol{x} \in \mathbb{R}^n$, $\boldsymbol{\lambda}$ and $\boldsymbol{\mu} \geq \mathbf{0}$. Taking the infimum on both sides then shows that

$$\inf_{\boldsymbol{x} \in \mathbb{R}^n} L_P(\boldsymbol{x}) \geq \inf_{\boldsymbol{x} \in \mathbb{R}^n} \{f(\boldsymbol{x}) + \boldsymbol{\lambda} \cdot (\mathbf{h}(\boldsymbol{x}) - \mathbf{c}) + \boldsymbol{\mu} \cdot \mathbf{g}(\boldsymbol{x})\} = L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}),$$

for every $\boldsymbol{\lambda}$ and $\boldsymbol{\mu} \geq \mathbf{0}$, and thus also $\inf_{\boldsymbol{x} \in \mathbb{R}^n} L_D(\boldsymbol{x}) \geq \sup_{\{(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\mu}}); \tilde{\boldsymbol{\mu}} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu})$. \square

Remark 18. Provided one is able to derive an expression for L_D , one is therefore in the very favorable setting of having to maximize a concave function subject to at most a simple non-negativity constraint. In addition, even if equality does not hold in (5.10), the solution of the dual problem does always provide a lower bound for that of the primal one. \diamond

Remark 19. Assuming there is an actual relation with the primal problem, the dual problem also provides an intuitive explanation for the relevance of the three conditions in Equations (5.2), (5.3) and (5.4). Since the definition of $L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\boldsymbol{x} \in \mathbb{R}^n} L(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ in Equation (5.9) is based on a free minimization of \boldsymbol{x} given $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$, the derivative with respect to \boldsymbol{x} should vanish. Similarly, as there is no restriction on $\boldsymbol{\lambda}$ in the dual problem (5.8), one might expect for the analogous condition to hold for $\boldsymbol{\lambda}$. Finally, with $\boldsymbol{\mu}$ being restricted to be non-negative, realizing the maximum of L_D with respect to $\boldsymbol{\mu}$ does not necessarily require for $\nabla_{\boldsymbol{\mu}} L(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{g}(\boldsymbol{x})$ to be zero, but should there be any positive entry $g_i(\boldsymbol{x})$, one could hope to further increase L by increasing μ_i .

It should be kept in mind though that, despite its intuitive appeal at first sight, this simple argument hides a critical point, namely that even though the inner minimization in \boldsymbol{x} is indeed a free one, the value of the minimizer will depend on the parameters $(\boldsymbol{\lambda}, \boldsymbol{\mu})$, i.e. the maximization of the dual function L_D with respect to $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ is in fact based on the function

$$L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) = L(\boldsymbol{x}(\boldsymbol{\lambda}, \boldsymbol{\mu}), \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

That it still makes sense requires a reasoning similarly to Remark 3, namely that despite the dependence of \boldsymbol{x} on the parameters, the optimality satisfied by \boldsymbol{x} ensures that L itself is in first order not affected by this dependence as the contributions $\frac{\partial L}{\partial \boldsymbol{x}} \frac{\partial \boldsymbol{x}}{\partial \boldsymbol{\lambda}}$ and $\frac{\partial L}{\partial \boldsymbol{x}} \frac{\partial \boldsymbol{x}}{\partial \boldsymbol{\mu}}$ arising in the derivatives with respect to the multipliers drop out by the optimality condition on \boldsymbol{x} .

The fact that a “blind” differentiation of the Lagrangian - i.e. ignoring any potential interplay between the various variables - is therefore still expected to deliver the correct results is clearly a major simplification, and likely one of the main reasons of its popularity in the engineering community. The same type of underlying “variational consistency” is, as in Remark 3, also fundamental to the ubiquitous changes of unknowns in thermodynamics¹. \diamond

Even though the lower bound in Lemma 2 can by itself be quite useful, the most desirable case is of course when both values actually do coincide. This motivates the following

¹The phasefield method owing much of its early success to its successful application in a thermodynamically based setting where such dependencies can legitimately be “ignored”, this is a point which unfortunately seems to be partially forgotten in the phasefield community.

Definition 4. (Saddle point)

A pair $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ is called a saddle point of the function $L(\mathbf{x}, \boldsymbol{\lambda})$ if one has

$$L(\bar{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) \leq L(\mathbf{x}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}). \quad (5.11)$$

for all $\mathbf{x} \in \mathbb{R}^n$ and $\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}$.

The next result is then a simple consequence of the Inequality (5.10):

Theorem 4. (thms. 4.8 and 4.9 [58])

If the Lagrangian $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ has a saddle point $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$, then $\bar{\mathbf{x}}$ is a solution of the primal problem, $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ is a solution of the dual problem, and one has

$$\max_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x}).$$

Conversely, assume that this relation holds true with finite values on both sides. Then, for every solution $\bar{\mathbf{x}}$ of the primal problem and every solution $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ of the dual problem, the point $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ is a saddle point of the Lagrangian.

Proof. If $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ is a saddle point of L , it follows from Equation (5.11) that one has

$$L_P(\bar{\mathbf{x}}) = \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L(\bar{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) \leq \inf_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) = L_D(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$$

and thus also

$$\inf_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x}) \leq L_P(\bar{\mathbf{x}}) \leq L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) \leq L_D(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) \leq \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}).$$

As the opposite inequality $\inf_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x}) \geq \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu})$ always holds as seen above, the inequalities necessarily hold as equalities,

$$\inf_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x}) = L_P(\bar{\mathbf{x}}) = L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) = L_D(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) = \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}). \quad (5.12)$$

Thus both the primal and dual problem do admit at least one solution $\bar{\mathbf{x}}$ and $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ and their values coincide, proving the first claim.

Conversely, if one has $\sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L_D(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathbb{R}^n} L_P(\mathbf{x})$ and $\bar{\mathbf{x}}$ resp. $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ are solutions to the primal resp. dual problem, one clearly has

$$L(\bar{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}); \boldsymbol{\mu} \geq \mathbf{0}\}} L(\bar{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = L_P(\bar{\mathbf{x}}) = L_D(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) = \inf_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) \leq L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$$

for all $\mathbf{x} \in \mathbb{R}^n$ and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ with $\boldsymbol{\mu} \geq \mathbf{0}$. Inserting $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ into the left-most expression and $\bar{\mathbf{x}}$ into the right-most one, it follows in addition that $L_P(\bar{\mathbf{x}}) = L_D(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}) = L(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$. \square

Remark 20. Note that this theorem does not assert the existence of a saddle point². Often, conditions assuring this are examined slightly more indirectly in terms of the **duality gap**, i.e. the difference between the inf and the sup obtained from the primal and dual problem³ or in relation with the Legendre-Fenchel transform introduced in the next section (see e.g. [26] and [14] for some links between the two approaches to duality). \diamond

²A fairly general existence theorem is e.g. given in [26].

³This of course amounts to the same as there being a saddle point is then clearly the same thing as this difference vanishing.

Nevertheless, it indicates one reason for the usefulness of the Lagrangian. Whereas the primal problem is by construction equivalent to the original minimization problem, this reformulation essentially has no practical use by itself. More precisely, if the constraints are satisfied for some \mathbf{x} , the values of the multipliers are, beyond the complementarity conditions for $\boldsymbol{\mu}$ essentially arbitrary. If the constraints are not satisfied, the inner supremum results in the value $+\infty$, which one does certainly not want to achieve in any practical minimization algorithm. It may thus not seem like a very reasonable idea to move the multipliers into the corresponding direction. In contrast, the statement implies that if there is a saddle point and one can find a solution $(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ of the dual problem, one has, at least, made the solution of the primal problem significantly simpler since it then reduces to a free minimization of $L(\mathbf{x}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ for the given values of the multipliers. In addition, provided this calculation can be justified through a sufficiently smooth dependence of the minimizer $\mathbf{x}(\boldsymbol{\lambda}, \boldsymbol{\mu})$ given some current estimate $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ in the definition of the dual function L_D in Equation (5.9), one has $\frac{\partial L}{\partial \mathbf{x}} = \mathbf{0}$ (as \mathbf{x} is a free minimizer for the given multipliers), and the **total** derivatives of L with respect to $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ therefore reduce to the partial ones despite this implicit dependence. It follows that

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}(\boldsymbol{\lambda}, \boldsymbol{\mu}), \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{h}(\mathbf{x}) - \mathbf{c} \quad \text{and} \quad \nabla_{\boldsymbol{\mu}} L(\mathbf{x}(\boldsymbol{\lambda}, \boldsymbol{\mu}), \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}).$$

This motivates for example various alternative descent/ascent methods in terms of the two sets of unknowns \mathbf{x} and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$, since, provided \mathbf{x} is at least an approximate minimizer for the given multipliers, adjusting them such as to increase L “as if” \mathbf{x} were fixed therefore does indeed make sense.

Remark 21. Even though often not primarily motivated by such considerations, the Lagrangian point of view also provides additional insight into e.g. some popular projection-based algorithms such as e.g. fractional-step algorithms in fluid dynamics and the return mapping algorithm in elasto-plasticity. In addition, the Lagrangian is also a very helpful ingredient for Newton-type schemes for constrained problems and the analysis of second-order necessary conditions for constrained optimization problems (see e.g. [58], [13] and [46] for a more detailed discussion). \diamond

Remark 22. It should be noted that the discussion above is for the most part completely independent of whether the underlying spaces are finite-dimensional or not. In particular, the definitions and lemmas can essentially be applied verbatim to the case when function spaces are involved, since e.g. the proofs of the concavity and the inequalities in Equation (5.10) in Lemma 2 as well as in Theorem 4 are all based purely on relations implied by the inf- and sup-operation (for Lemma 2 combined with the concavity of linear operators), which are completely independent of any particularly favorable properties of \mathbf{R}^n and neither rely on any particular structure of the sets involved in the inf-sup-operations. For the corresponding definitions and proofs in this more general setting, a classical and very readable reference is [26]. \diamond

Remark 23. The dual function $L_D(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is an example of a more general kind of function frequently arising in parametric optimization involving a function $\varphi(\mathbf{x}, \mathbf{u})$ of the variable \mathbf{x} and the “parameter” \mathbf{u} . Based on $\varphi(\mathbf{x}, \mathbf{u})$, one can introduce the (optimal) value function $v(\mathbf{u}) := \inf_{\mathbf{x} \in \mathbf{R}^n} \varphi(\mathbf{x}, \mathbf{u})$ corresponding to a minimizer (if any) of φ given the value of the parameter \mathbf{u} . More generally, one might also restrict the domain of \mathbf{x} in the minimization process just to subsets the underlying space, i.e. by setting $v(\mathbf{u}) := \inf_{\mathbf{x} \in X} \varphi(\mathbf{x}, \mathbf{u})$, or even to subset $X(\mathbf{u})$ depending themselves on the parameter.

It is clear that the study of the behavior of $v(\mathbf{u})$ in terms of \mathbf{u} (e.g. regularity or differentiability properties), has an inherent interest beyond the Lagrangian setting above, but is in this particular setting also very instructive in terms of an additional interpretation of the Lagrange-multipliers as “sensitivities” of the objective functions with respect to the constraints. For further discussions in the finite-dimensional setting, the reader is e.g. referred to [46], [13], [58] and [27]. For the infinite-dimensional setting in function spaces, the reader may consult [26] or for a more detailed but also more technical discussion in a quite general setting [14]. \diamond

5.2 The Legendre-Transform

An a priori somewhat different but in fact closely related approach to duality theory is provided through the **Legendre-transform**, which is in particular heavily relied upon in the thermodynamic setting. In its most basic form, the Legendre transform (or **conjugate**) of a function $f : \mathbb{R} \rightarrow (-\infty, +\infty]$ is a second function $f^* : \mathbb{R} \rightarrow (-\infty, +\infty]$ defined as⁴

$$f^*(x^*) := \sup_{x \in \mathbb{R}} \{x^*x - f(x)\}, \quad (5.13)$$

i.e. the function f^* which, for an arbitrary “slope” x^* , assign to f the maximal difference between the line $y(x) = x^*x$ and the graph of f . One can reiterate this procedure on $f^*(x^*)$ by defining the **bi-conjugate (bi-dual)** function

$$f^{**}(x) := \sup_{x^* \in \mathbb{R}} \{x^*x - f^*(x^*)\}, \quad (5.14)$$

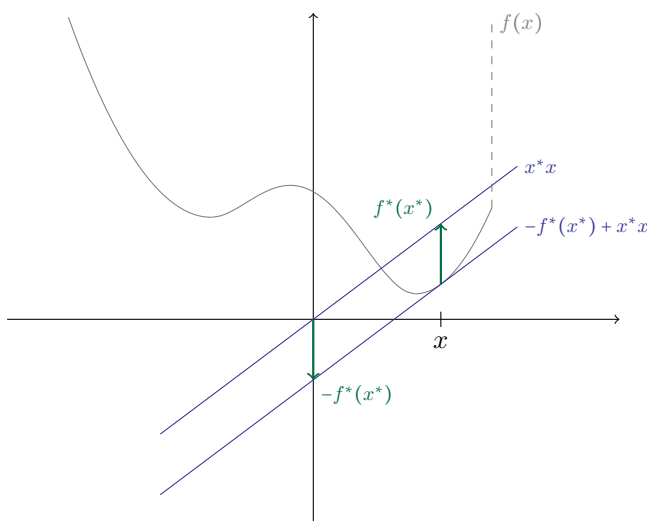


Figure 5.1: Illustration of the Legendre transform

The motivation for these definition is illustrated in Figure 5.1: $f^*(x^*)$ being the largest difference between $y(x) = x^*x$ and $y = f(x)$, the line $y(x) = -f^*(x^*) + x^*x$ necessarily passes below the graph (in the sense $y(x) \leq f(x)$ for all x). In addition, for the particular point chosen in Figure 5.1, this line is actually exact at x , meaning that $f(x) = -f^*(x^*) + x^*x$, with a slope given precisely by $f'(x)$ as it would otherwise cut the graph. The Legendre transform can therefore be interpreted as implicitly constructing tangent lines to the graph of a function for each given slope, at least provided this is possible without cutting the graph.

In addition, comparing the graphs of $-f^*(x^*) + x^*x$ for various values of x^* in Figure 5.2, this translated line actually is the one achieving the maximum value for the given value of x as in the definition of f^{**} since, by the very definition of f^* , all other lines $-f^*(x^*) + x^*x$ necessarily

⁴Note that the value $+\infty$ is explicitly admissible here for f . One could in principle also think about admitting $-\infty$, but excluding this value on the one hand simplifies certain statements and on the other hand does not really eliminate any interesting functions, since a convex function which is $-\infty$ at some point can at most be finite at a single point before jumping to $+\infty$ as convexity enforces an infinite slope (similar to the jump to infinity in the right of Figure 5.1).

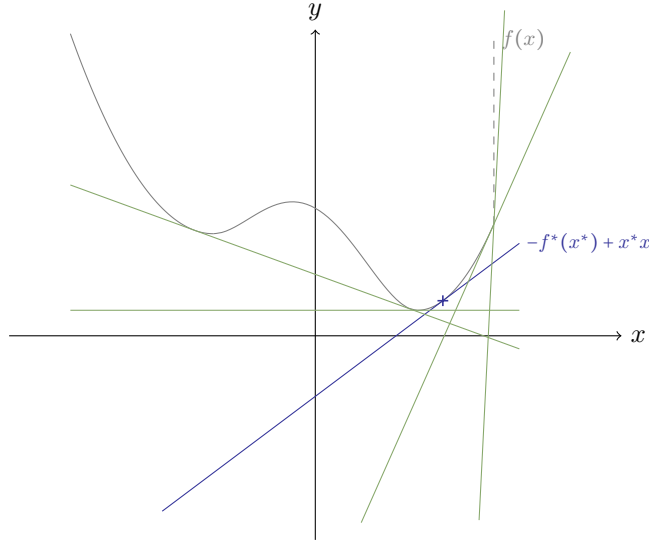


Figure 5.2: Illustration of the inverse Legendre transform

pass below the point $(x, f(x))$. This shows that that one can, at least in certain situations⁵, recover the value $f(x)$ through $f^*(x^*)$.

As it turns out, the arguments underlying the theory for this type of transform have virtually nothing to do with f being defined on \mathbb{R} . More generally, the same procedure can be applied for a multivariate function $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ by replacing the line $(x, x^* \cdot x)$ with the hyperplane $(\mathbf{x}, \mathbf{x}^* \cdot \mathbf{x})$ in \mathbb{R}^{n+1} , $f^*(\mathbf{x}^*) = \sup_{\mathbf{x} \in \mathbb{R}^n} \{\mathbf{x}^* \cdot \mathbf{x} - f(\mathbf{x})\}$, or, if X is e.g. a Banach space⁶ with dual X^* , and $f : X \rightarrow (-\infty, +\infty]$ a given functional, one can similarly define

$$f^*(x^*) := \sup_{x \in X} \{ \langle x^*, x \rangle - f(x) \}, \quad (5.15)$$

where $\langle x^*, x \rangle$ denotes the natural pairing between X and X^* . Note in particular the important relation

$$x^* \cdot x \leq f(x) + f^*(x^*), \quad (5.16)$$

valid for arbitrary $x \in X$ and $x^* \in X^*$. One can reiterate the procedure - but by again using X instead of $(X^*)^* = X^{**}$ - by defining the **bi-dual function**⁷

$$f^{**}(x) := \sup_{x^* \in X^*} \{ \langle x^*, x \rangle - f^*(x^*) \}, \quad (5.17)$$

whose primary motivation stems from the following remarkable property:

Theorem 5. (Fenchel-Moreau, thm. 1.11 [16])

Assume that $f : X \rightarrow (-\infty, +\infty]$ is convex, lower semicontinuous and not identically equal to $+\infty$. Then $f^{**} = f$.

⁵This will be made more precise below. As might be guessed from Figure 5.1 and the discussion thus far, it turns out that the crucial condition is that there is in fact a line passing below the graph of f whose value coincides with $f(x)$, i.e. a so-called subgradient (see in particular Props. 2 and 3 below).

⁶Note that the norm of X is not directly used in any of the results below, i.e. one can consider, as in e.g. [26] or [14], even more general spaces.

⁷The distinction between X^{**} and X is of course irrelevant if X is reflexive, i.e. in particular in \mathbb{R}^n or the Hilbert-space setting. It is only when the inclusion of X in X^{**} is strict that this restriction to X becomes important, as one is primarily interested in the relation between f^{**} and f , the latter one a priori only being defined on X .

Even if the assumptions above do not hold, one always has the following fundamental estimate:

Lemma 3. *Let $f : X \rightarrow (-\infty, +\infty]$ be arbitrary. Then $f^{**}(x) \leq f(x)$.*

Proof. This is a direct consequence of taking the supremum over the Equality (5.16) in the form $\langle x^*, x \rangle - f^*(x^*) \leq f(x) \quad \forall x \in X, x^* \in X^*$. \square

Remark 24. Before stating any further properties leading up to and/or derived from the Fenchel-Moreau-theorem, it is important to note that, besides the standard mathematical definition above, there are, in particular in the physical setting, various alternative definitions of Legendre-type transformations. A very similar one to (5.15) (again primarily useful if f is convex, see below) is given by

$$\hat{f}(x^*) = \inf_{x \in X} \{f(x) - \langle x^*, x \rangle\} = -\sup_{x \in X} \{\langle x^*, x \rangle - f(x)\} = -f^*(x^*) \quad (5.18)$$

enjoying similar properties, but with \hat{f} now being a concave function (as in infimum of linear and thus concave functions in x^*). Under the same assumptions on f as in Thm. 5, the inversion formula (5.17) above can be rewritten as

$$f(x) = f^{**}(x) = \sup_{x^* \in X^*} \{\langle x, x^* \rangle - f^*(x^*)\} = \sup_{x^* \in X^*} \{\langle x^*, x \rangle + \hat{f}(x^*)\}. \quad (5.19)$$

If f itself is concave instead of convex, a more useful definition is obtained by setting e.g.

$$\tilde{f}(x^*) = \sup_{x \in X} \{\langle x^*, x \rangle + f(x)\} = \sup_{x \in X} \{\langle x^*, x \rangle - (-f)(x)\} = (-f)^*(x^*), \quad (5.20)$$

leading again to a convex function (note that $(-f)^* \neq -f^*$, which would be concave!). Again based on the conversion formula (5.17), under the appropriate assumptions, one has

$$f(x) = -((-f)^*)^* = -\sup_{x^* \in X^*} \{\langle x^*, x \rangle - (-f)^*(x^*)\} = \inf_{x^* \in X^*} \{(-f)^*(x^*) - \langle x^*, x \rangle\}. \quad (5.21)$$

Similar modifications are clearly also possible when reverting the signs of e.g. x^* in (5.15) or (5.18) as this just applies a change in sign of the respective arguments. \diamond

At first sight, the definitions above differ from the usual procedure in the physical literature (here for simplicity again in the one-dimensional setting) of introducing the variable x^* by setting $x^* := \frac{df}{dx}$ and then "defining" the Legendre transform of f e.g. as

$$f^*(x^*) := x^*x - f(x), \quad (5.22)$$

where the sole dependence of g on x^* is justified by remarking that $\frac{\partial}{\partial x}(x^*x - f(x)) = x^* - \frac{df}{dx} = 0$ as $x^* = \frac{df}{dx}$. An "inversion formula" is then recovered by simply rearranging Equation (5.22) to $f(x) = x^*x - f^*(x^*)$, where, by a simple differentiation of (5.22), x in addition satisfies $x = \frac{df^*}{dx^*}$. This is a somewhat tricky argument though, as x^* , if **defined** as $\frac{df}{dx}$, is clearly a function of x , $x^* = x^*(x)$, and f^* therefore in fact still an explicit function of x instead of x^* , $f^* = f^*(x)$, meaning that the argument above is only valid if one chooses to "forget" this (explicit) dependence.

The definition (5.15) on the one hand avoids this pitfall as f^* is evidently truly independent of x as the original variable is eliminated through the supremum operation, and is, on the other hand, more general as it does not require any differentiability of f .

If f is strictly convex and smooth, one can easily establish a link between the two approaches though. In fact, under this assumption, the supremum in (5.15) will be taking in a single point x , which is characterized through the Euler-Lagrange equation $\frac{\partial}{\partial x}(x^*x - f(x)) = x^* - \frac{\partial f}{\partial x}(x) = \mathbf{0}$, i.e. the relation $x^* = \frac{\partial f}{\partial x}$ above. The subtle but important difference in interpretation here is that this does not define x^* - which is given - as a function of x , but rather the point x realizing the supremum as an implicit function of x^* . A more proper way of writing (5.22) would therefore be as

$$f^*(x^*) = x^*x(x^*) - f(x(x^*)),$$

which is then obviously a function of x^* alone. Assuming the dependence of x on x^* is in addition differentiable, g satisfies

$$\frac{df^*}{dx^*} = x(x^*) + \left(x^* - \frac{df}{dx}(x(x^*))\right) \cdot \frac{dx}{dx^*}.$$

As the last term drops out due to $x(x^*)$ satisfying $x^* = \frac{df}{dx}(x(x^*))$, one thus recovers the previous relation $x(x^*) = \frac{\partial f^*}{\partial x^*}$, now as an explicit relation in terms of f^* .

Similarly, assuming that f^* is also strictly convex and smooth, the supremum in (5.17) will also be achieved in a single point x , whose Euler-Lagrange equation is given by (note x is again just an arbitrary **given** slope here) $\frac{\partial}{\partial x^*}(x^*x - f^*(x^*)) = x - \frac{\partial f^*}{\partial x^*}(x^*) = 0$, which then defines x^* as an implicit function of x . Using this and the fact that under the given assumptions $f^{**} = f$, one has $f(x) = f^{**}(x) = x^*(x)x - f^*(x^*(x))$, and, if the mapping $x \mapsto x^*(x)$ is actually differentiable, $\frac{df}{dx}(x) = x^*(x) + \left(x - \frac{df^*}{dx^*}(x^*(x))\right) \frac{dx^*}{dx}$. The last term again drops out by the definition of $x^*(x)$, i.e. one recovers the formula $x^*(x) = \frac{df}{dx}(x)$.

This type of differential relation can be substantially generalized to situations where f involving less smoothness⁸ This requires replacing the classical derivative of f with an appropriate generalized notion of differentiability, which, in the convex (resp. concave) setting is given by the following

Definition 5. (Subdifferential)

Let $f : X \rightarrow (-\infty, +\infty]$. f is said to be **subdifferentiable** (in the sense of convex analysis) at a point x in X if $f(x)$ is finite and there is some x^* in X^* such that the relation

$$f(y) \geq f(x) + \langle x^*, y - x \rangle \tag{5.23}$$

holds for all $y \in X$, i.e. if the hyperplane $(x, \langle x^*, x \rangle)$ in $X \times \mathbb{R}$ with slope x^* passing through the point $(x, f(x))$ lies below the graph of f . The set of all (if any) such slopes at a given point x is called the **(convex) subdifferential** ∂f of f at x . If $f(x)$ is not finite, ∂f is defined to be empty, i.e.

$$\partial f(x) = \begin{cases} \emptyset & \text{if } f(x) \text{ is not finite,} \\ \{x^* \in X^* : f(y) \geq f(x) + \langle x^*, y - x \rangle \forall y \in X\} & \text{else,} \end{cases}$$

and f is thus subdifferentiable at x iff $\partial f(x) \neq \emptyset$.

Similarly, the **concave subdifferential** $\partial^\cap f$ of f is obtained by reversing the sign of the inequality in Equation (5.23), and therefore corresponds to the slopes of all hyperplanes lying above the graph of f and passing through $(x, f(x))$.

For convex functions, an important relation with the more basic notion of Gâteaux-differentiability is given by the following

⁸Even though it is well-known that the convexity itself actually does imply a certain degree of regularity, at least within the interior of the domain of f .

Proposition 1. (Subdifferentiability and Gâteaux-differentiability, prop. I.5.3 [26])

Let $f : X \rightarrow (-\infty, +\infty]$ be convex. If f is Gâteaux-differentiable at x , it is subdifferentiable at x . Conversely, if f is continuous and finite at x and has only one subgradient, then f is Gâteaux-differentiable at x . In both cases, the relation $\partial f(x) = \{f'(x)\}$ holds.

With this definition at hand, the generalization of the relation $x^* = \frac{\partial f}{\partial x}$ above is contained in the following

Proposition 2. (see e.g. prop. 5.1. [26])

Let f be a function of $X \rightarrow (-\infty, +\infty]$ and f^* its conjugate, and further assume that f is not identically equal to $+\infty$. Then $x^* \in \partial f(x)$ iff

$$f(x) + f^*(x^*) = \langle x^*, x \rangle. \quad (5.24)$$

Proof. Assume that (5.24) holds. The first - somewhat technical - observation is that this implies that both $f(x)$ and $f^*(x^*)$ are finite. In fact, as $\langle x^*, x \rangle$ is always finite⁹ and f has by assumption at least one point where it is finite, f^* is never equal to $-\infty$. The sum of $f(x)$ and $f^*(x^*)$ being finite and opposite infinities not being possible, both need to be finite. If the relation (5.24) holds, the hyperplane $(x, -f^*(x^*) + \langle x^*, y \rangle)$ passes through $(x, f(x))$ (i.e. is exact at x). As Equation (5.24) together with the inequality (5.16) imply that x in fact realizes the supremum in the definition of f^* , it holds that, for any $y \in X$, $f^*(x^*) = \langle x^*, x \rangle - f(x) \geq \langle x^*, y \rangle - f(y)$, and thus by rearrangement both $f(y) \geq -f^*(x^*) + \langle x^*, y \rangle \quad \forall y \in X$ (i.e. the graph of f lies entirely above this hyperplane) and $f(y) \geq f(x) + \langle x^*, y - x \rangle \quad \forall y \in X$ (i.e. the defining inequality (5.23) for a subgradient holds).

Conversely, if $x^* \in \partial f(x)$, by the same inequality (5.24), one has $\langle x^*, x \rangle - f(x) \geq \langle x^*, y \rangle - f(y)$ for all $y \in X$. Taking the supremum over y shows that

$$\langle x^*, x \rangle - f(x) = \sup_{y \in X} \{ \langle x^*, y \rangle - f(y) \} = f^*(x^*).$$

□

Despite the suggestive arrangement of the duality relation (5.24), it is not quite symmetric in f and f^* , as f^* is the transform of f , but f need not be the one of f^* unless the inversion formula (5.17) holds. In fact, applying Proposition 2 starting from f^* instead of f a priori only shows that

$$x \in \partial f^*(x^*) \Leftrightarrow f^*(x^*) + f^{**}(x) = \langle x^*, x \rangle. \quad (5.25)$$

Whether or not a “symmetric” version of Proposition 2 holds thus clearly depends on the relation between f and f^{**} . This is clarified in the following

Proposition 3. (Cor. 5.2 [26], Prop. 2.118 [14])

Let $f : X \rightarrow (-\infty, +\infty]$ be a (possibly non-convex) function. Then the following holds:

1. If $x^* \in \partial f(x)$, then $x \in \partial f^*(x^*)$.
2. If f is subdifferentiable at x , then $f^{**}(x) = f(x)$.
3. If $f^{**}(x) = f(x)$, then $\partial f(x) = \partial f^{**}(x)$ (which is empty if both are equal to $\pm\infty$) and the stronger statement

$$x^* \in \partial f(x) \Leftrightarrow f(x) + f^*(x^*) = \langle x^*, x \rangle \Leftrightarrow x \in \partial f^*(x^*), \quad (5.26)$$

i.e. the statement of 5.24 together with its “dual” (instead of 1), holds. In addition, if the common value $f^{**}(x) = f(x)$ is finite, the variational characterizations

$$\partial f(x) = \operatorname{argmax}_{x^* \in X^*} \{ \langle x^*, x \rangle - f^*(x^*) \}. \quad (5.27)$$

⁹By assumption, $x^* \in X^*$ and is thus a **continuous** linear functional.

and

$$\partial f^*(x^*) = \operatorname{argmax}_{x \in X} \{\langle x^*, x \rangle - f(x)\}. \quad (5.28)$$

hold for the subgradients of f and f^* .

Proof. The proof is mostly a variation of the arguments in the proof of Proposition 2 combined with Lemma 3.

As shown above, $x^* \in \partial f(x)$ iff Equation (5.24) holds, in which case $f(x^*)$ is necessarily finite, and the hyperplane $-f(x) + \langle x^*, x \rangle$ is exact for f^* at x^* . That the graph of $f^*(x^*)$ lies above this hyperplane now follows from Equation (5.16) since $\langle x^*, x \rangle - f(x) \leq f^*(x^*)$ holds for all $x \in X^*$. As to the second point, if $\partial f \neq \emptyset$, by Prop. 2, there exists some $x^* \in X^*$ such that

$$f(x) = \langle x^*, x \rangle - f^*(x^*) \leq \sup_{x \in X} \{\langle x^*, x \rangle - f^*(x^*)\} = f^{**}(x),$$

whereas the reverse equality always holds by Lemma 3. Finally, as $x^* \in \partial f(x)$ iff $f(x) = \langle x^*, x \rangle - f^*(x^*)$ and $x^* \in \partial f^{**}(x)$ iff $f^{**}(x) = \langle x^*, x \rangle - f^*(x^*)$, $f(x) = f^{**}(x)$ finite implies that both conditions necessarily hold at the same time. Due to the equality of $f(x)$ and $f^{**}(x)$ as well as $\partial f(x)$ and $\partial f^{**}(x)$, Equation (5.26) follows directly from Proposition 2. As already seen in the proof of this proposition, $x^* \in \partial f(x)$ iff x realizes the supremum in the definition of f^* , which, by the symmetry in the case above, happens iff $x \in \partial f^*(x)$, which in turn happens iff x^* realizes the supremum in the definition of $f^{**}(x) = f(x)$. \square

Remark 25. Note that, more generally, if, for some $x \in X$, the value $f^{**}(x)$ is finite, then the variational characterization

$$\partial f^{**}(x) = \operatorname{argmax}_{x^* \in X^*} \{\langle x^*, x \rangle - f^*(x^*)\}$$

holds for the (potentially empty) subgradient of f^{**} at x based on considering the duality between f^{**} and $f^{***} = f^*$ (which always holds)¹⁰. If $f^{**}(x) \neq f(x)$, there is a priori no need for their subgradients to coincide, and the implication $x \in \partial f^*(x^*) \Rightarrow x^* \in \partial f(x)$ need **not** hold as, by Equation (5.25), the hyperplane $(x, -f^*(x^*) + \langle x^*, x \rangle)$ is exact at $f^{**}(x)$, but not at $f(x)$. In fact, as f^{**} is the pointwise supremum of all affine functions lying below the graph of f , the subgradient $\partial f(x)$ has to be empty at such points. Otherwise any subgradient would be based by definition on such an affine function which in addition passes through $(x, f(x))$, therefore precluding $f^{**}(x) < f(x)$. \diamond

A particularly important example for a non-smooth situation where full duality holds in the optimization setting is given by the **indicator function** of a nonempty closed convex subset K ,

$$I_K(x) = \begin{cases} 0 & x \in K, \\ +\infty & \text{else.} \end{cases}$$

This can be shown to be a convex, lower semicontinuous function due to the assumptions on K . As the supremum defining the Legendre-transform is always equal to $-\infty$ if $x \notin K$, it is easy to see that

$$I_K^*(x^*) = \sup_{x \in K} \{\langle x^*, x \rangle\} = \Pi_K(x^*),$$

where, for any subset S of \mathbb{R}^n , the function $\Pi_S(x^*) := \sup_{x \in S} \{\langle x^*, x \rangle\}$ is the so-called **support function** of the set S . By the basic properties of the Legendre-transform, Π_K is again a convex lower semicontinuous function if K is closed and convex, and, from Thm. 5, $\Pi_K^*(x) = I_K(x)$. The

¹⁰Note that the proof of this statement in [14] is slightly misleading, as applying their eq. 2.229 with f^* instead of f^{**} replacing f would only imply the equivalences (5.25) instead of the equivalence with $x^* \in \partial f^{**}(x)$.

subgradient of I_K is empty whenever $x \notin K$ (as the value of I_K is not finite there), and, directly applying the condition (5.23) as the set of all $x^* \in X^*$ such that

$$I_K(y) \geq I_K(x) + \langle x^*, y - x \rangle \stackrel{I_K(x)=0}{=} \langle x^*, y - x \rangle.$$

Since this inequality is trivially satisfied for $y \notin K$ as $I_K(y) = +\infty$, the only relevant condition to check is what happens if $y \in K$ (and thus $I_K(y) = 0$), and the subgradient is thus given by all x^* such that $0 \geq \langle x^*, y - x \rangle$ for all $y \in K$. This situation can be summarized by defining the **normal cone** to K at an arbitrary x as

$$N_K(x) = \begin{cases} \{x^* \in X^* : \langle x^*, x - y \rangle \geq 0 \ \forall y \in K\} & x \in K, \\ \emptyset & \text{else,} \end{cases} \quad (5.29)$$

with which one has $\partial I_K(x) = N_K(x)$. Based on Prop. 2, any of the subgradients of $\Pi_K(x^*)$ is characterized by the equality $I_K(x) + \Pi_K(x^*) = \langle x^*, x \rangle$, i.e. $\sup_{y \in K} \{\langle x^*, y \rangle\} = \langle x^*, x \rangle - I_K(x)$. If $x \notin K$, the right-hand side equals $-\infty$, and as K is nonempty this equality cannot hold. All potential subgradients thus have to lie in K . If $x \in K$, the indicator-function drops out, and the subgradient of Π_K at x^* is given by the set of all x such that $\sup_{y \in K} \{\langle x^*, y \rangle\} = \langle x^*, x \rangle$, i.e. in accordance with Equation (5.28) as $\partial \Pi_K(x^*) = \operatorname{argmax}_{x \in K} \{\langle x^*, x \rangle\}$ and therefore the vector(s) in K furthest away from the origin “along” x^* .

Remark 26. Note that the first conclusion is in principle another very efficient way of obtaining the first-order necessary (and sufficient) conditions for convex constrained optimization problems and that the supremum in the definition of the primal function $L_P(\mathbf{x}) = \sup_{\{(\boldsymbol{\lambda}, \boldsymbol{\mu}) : \boldsymbol{\mu} \geq \mathbf{0}\}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ in Equation (5.6) has precisely the same effect as adding the indicator function for the admissible set to f .

In fact, it is easy to see that that the necessary and sufficient condition for a convex minimization problem is given by $0 \in \partial f(x)$, since this means that there is a horizontal “plane” passing everywhere below (in the sense of \leq) the graph and touching it at the point x , and x therefore has to be a minimizer (see the horizontal line in Figure 5.1). If f is smooth on K , one has¹¹ $\partial(f + I_K)(x) = \{f'(x)\} + \partial I_K$, and requiring that 0 lie in this set therefore shows that for any minimizer one must have $-f'(x) \in N_K(x)$.

Nevertheless, the Lagrangian formulation, while more susceptible to failure unless a constraint qualification condition is satisfied, has the major advantage of directly delivering a “good guess” of what the normal cone looks like based on an algebraic characterization of the constraint set, whereas this question is left open in the characterization of the normal cone in Equation (5.29)¹². Some further links between the two approaches will be summarized in the next section. \diamond

¹¹Note that the equality $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$ is not generally true but holds under some mild regularity condition of f , see e.g. [26] or [20].

¹²It should be kept in mind that, similar to Remark 6, the normal cone is the actually relevant set, which can often - but not always - be characterized using the gradient of the functions defining the constraint.

Part III

Applications

The phasefield method is primarily an approach used for the modeling and simulation of phase transformations within the material sciences. In particular, it has been and continues to be heavily employed as a tool to gain a deeper understanding of the effects of different process conditions on the resulting microstructure during solidification processes. Additional - somewhat more recent - applications include the study of the evolution of purely solid microstructures when e.g. subjected to varying external loads or, more simply, during an aging process.

While many of the earlier phasefield models were already able to provide a qualitative understanding of the evolution of microstructures, it was soon realized that there are some serious practical difficulties when trying to obtain more quantitative insights. These are primarily due to the very idea underlying the phasefield method in combination with limited computational resources.

On the one hand, the lengthscale ϵ_p associated with a physical interface region between two different phases is usually several orders of magnitude smaller than the lengthscale L associated with the microstructure itself. Even though this makes the volume of the interfaces essentially negligible, these are nevertheless regions associated with a high energy density due to the mismatch in the atomic arrangement between the materials on both sides. This energy contribution can therefore normally not be neglected but can, at the scale L of the microstructure, be in a good approximation be associated with a surface energy density $\gamma^{\alpha\beta}$ on the “almost” two-dimensional interface. The major advantage of this approximation is that the total energy (resp. entropy) of the microstructure can then be described in terms of volumetric contributions from the various pure phases in terms of their bulk-properties and the effective surface energy densities $\gamma^{\alpha\beta}$ hiding the highly complex physics within the true interfaces. In addition, being an effective macroscopic property, these surface energy densities may in particular be experimentally accessible and can therefore serve as input for numerical models for the evolution of the microstructure.

On the other hand, numerically modeling the evolution of a problem defined by moving (sharp) interfaces is quite challenging. The core of the phasefield method when applied to such a situation lies in “partially undoing” this sharp interface limit by reintroducing a small but finite transition region of width $\sim \epsilon$ between the various phases. Given that the phasefield functional can be based on the **measured** effective (macroscopic) surface energy densities $\gamma^{\alpha\beta}$, one has the great advantage that it is neither necessary for this artificial length scale ϵ to match with the actual width ϵ_p of the physical interface, nor to include a complex microscopic model for this transition region.

Nevertheless, ϵ can of course not be chosen arbitrarily large as it needs to be sufficiently small in order for the phasefield functional to provide an accurate approximation of the energy associated with the sharp interface setting. As the convergence of the “pure” phasefield functional to the surface energy basically relies on the phasefield profile within the interface converging to the classical one-dimensional steady-state profiles (see Section 6.2.2), in combination with additional driving forces arising in the coupled models¹³, this essentially imposes at least three restrictions. The first one is purely geometrical and requires for the interface width to be significantly smaller than the (smallest) radius of curvature. Otherwise, the “tangential” derivatives of the phasefield (w.r.t. to an assumed sharp interface surface lying e.g. at the 0.5 isoline) may become competitive with the “normal” ones as one moves away from this line, thus leading to a distorted profile. The second and third are somewhat related and depend upon the additional energetic contributions and their associated driving forces. On the one hand, even in the absence of any curvature, it is usually necessary for the interface width to be significantly smaller than the total length scale of the problem as otherwise the contribution of the - a priori to a degree arbitrary - interpolation of the given bulk energy densities within the interface region can make up a significant part of the total energy and thus lead to a large distortion of the energetics unless the

¹³The same issue arises can even arise in the absence of additional energetic contributions in the pure phasefield case due to additional constraints such as a volume constraint.

interpolated energy happens to coincide very closely with the “true” energy.

On the other hand, it is also necessary for the strength of the pure phasefield terms (i.e. the ones arising from a and w) to be locally significantly larger than these additional driving forces in order to avoid an excessive deformation of the interface profile.

Which of the three constraints on ϵ is more restrictive clearly depends upon the particular problem. The most obvious idea for avoiding any such issues is of course to simply choose a smaller interface width ϵ as compared to the size of the microstructure. It is clear that if one does so using a given grid-spacing, a reduction in ϵ will entail higher numerical errors in resolving the transition of the phasefield within the increasingly narrow interfaces. Even though this procedure is free of any computational cost (and will actually usually even decrease it) as it amounts to simply reducing a parameter and will improve the approximation at a continuous level, the numerical difficulties will at some point overshadow this improvement. Once this point is reached, the only viable alternative is then an increase of the resolution, which, unlike the modification of ϵ , will always entail an increase in the computational effort. Finding the “sweet spot” for a given problem, i.e. the point where the most accurate results are - due to a good balance between numerical and continuous (in terms of ϵ) approximation errors - obtained for a given cost is a fairly difficult problem, which also depends heavily on the interplay of various practical factors. Some of these as well as their interplay with some of the numerical difficulties will be discussed in Chapter 6, with a strong focus on an obstacle-potential based setting.

In contrast to these primarily numerical factors, the way any additional energy contributions $f(\phi, \mathbf{c}, T, \dots)$ affect the precision of the results for a given ϵ can often be influenced quite heavily by the particular manner they are modeled within the interface region. Since easing the restriction on the numerical interface width can lead to a significant decrease of the required resolution and therefore the necessary computational effort, this has been the subject of a fairly extensive research effort over the past decades and will be the primary focus of Chapter 7.

Chapter 6

Pure Phasefield Problems

The decisive role of the phasefield method within a modeling context is its ability to represent effects associated with surface energies (or entropies). Two very useful properties in this respect are that (at least when approaching the sharp-interface limit) this ability has relatively little to do with the presence of additional energetic contributions provided these are well-behaved as $\epsilon \rightarrow 0$, and that these additional influences can be included by simply “adding” them to the essential gradient energy density and bulk-potential contributions. While this statement needs to be put into perspective as it is in practice often necessary to use artificially enlarged interfaces due to limited computational resources¹, it explains the basic additive structure with respect to the driving forces underlying the common phasefield modeling approach. In addition, it highlights the crucial importance of these “standard” phasefield contributions in general.

Before discussing the more complex multiphysics problems alluded to in Chapter 3, it is therefore instructive to first consider the simpler setting in which the phasefield variable is the only unknown. On the one hand, this will serve to introduce some more details on the particular phasefield model used in this thesis. On the other hand, while coupling the phasefield and other fields - even those governed by more “classical” equations which are better understood at both a theoretical and numerical level - can often lead to additional challenges by itself, there is also a number of relatively generic challenges associated with phasefield-type problems, in particular from a numerical point of view. Two of these - related with the bound-constraints in the obstacle case - will be discussed in more detail in this chapter.

Remark 27. The “pure” phasefield-setting in this chapter is chosen primarily due to its simplicity and since there is no point in adding any additional complexity through couplings for the discussion below. Nevertheless, the minimization of surface energy is an interesting phenomenon by itself, and there is in fact a number of relevant applications which can be modeled using the phasefield variable alone. Two particularly interesting ones are the equilibrium shape of droplets on substrates or fibers, for which the interested reader is referred to e.g. [9] and [1]. \diamond

Remark 28. A sizeable part of the research focus in the phasefield community in the last decades has been oriented towards improving the accuracy of phasefield methods despite the practical need for artificially enlarged interfaces (often in terms of thin interface limits), in particular for coupled multiphysics problems. This can often be achieved through an improved modeling of the precise form of additional energy contributions within the interfacial regions (two such examples will be considered in Sections 7.1 and 7.2). Even though this can lead to a quite complex description, this added complexity has very little direct influence on the discussions in this chapter. \diamond

¹Using a small but finite interface width will generally lead to deviations from the desired sharp-interface limits. These are often related with the problem of **excess energies** and/or high **Cahn-numbers** (the ratio of the interface width to the “radius” or, in the non-circular case, a measure thereof).

Remark 29. There are also applications where the phasefield method is not primarily used due to its ability to capture effects associated with surface energies, but instead as a purely numerical tool with the main purpose of replacing a physical problem involving sharp interfaces (whether moving or not) with a diffuse approximation thereof. This includes in particular models where the surface-minimizing property of the standard phasefield model is eliminated through an appropriately chosen “counter-term” with the purpose of obtaining a method implicitly tracking sharp interface motions (see e.g. [72], [67] and the references therein for some background and applications). Even though this can be of obvious interest from a numerical point of view, as it in principle allows for working on fixed grids (in particular simple Cartesian ones) even in the presence of complex and potentially moving interface, such problems will not be considered here. \diamond

Section 6.1 will first provide a quick outline of the basic variational framework in the multiphase case, corresponding essentially to a simplified version of the one in [52] underlying most the work in this thesis to a pure phasefield setting. Before returning to this more complex case in Section 6.3, the discussion in Section 6.2 will again focus on the simplest possible setting consisting of a reduced (i.e. expressed solely in terms of a single phasefield ϕ) two-phase version of this general setting. After introducing some standard background and results, the main focus of this section will be a relatively detailed analysis of the impact of the 0-1-bounds on the phasefield values in the discrete case as one of the central “ingredients” of the basic phasefield model in the obstacle case. More precisely, a discrete equivalent of the basic analytical one-dimensional phasefield profile (corresponding to an undisturbed flat interface) and the associated energetics will be derived for the case of the obstacle potential. On the one hand, this allows for a very simple interpretation of some commonly observed effects (such as the discrete gradient energy always being larger than the continuous one). On the other hand, the resulting expressions can conveniently be expanded in terms of the discretization parameter Δx , from which a number of interesting facts can simply be “read off”. In particular, even though the discrete interface width is only first-order accurate (and always more narrow than the continuous one), both the energetics and the discrete profile itself are second-order convergent.

Some additional issues related to the multiphase setting will then be discussed in Section 6.3. Subsection 6.3.1 will first recall the basic equations to be fulfilled and different choices for the dynamics. This is followed in Subsection 6.3.3 by a discussion of some numerical and algorithmical aspects arising when dealing with multiphase problems.

Remark 30. Parts of this chapter consist of elementary background for the phasefield equation. It was nevertheless chosen to introduce this background here instead of the actual background part since it is on the one hand mostly directly relevant for the discussion of the main points and on the other hand, explicitly being based on the phasefield equation itself, considerably more specific than the relatively general and consequently abstract previous considerations.

The intention is also not to provide a full discussion of general phasefield problems (which, in particular in the multiphase case, would be a very complex and difficult task), but simply to recall some elementary facts before focusing on some particular but quite relevant topics encountered when using an obstacle potential. \diamond

6.1 The Basic Phasefield Functional

The key idea underlying the phasefield method for multiphase problems is essentially the same as the one already outlined for the simpler two-phase setting in section 2 and consist in the use of a vectorial **order parameter** $\phi = (\phi^\alpha)_{1 \leq \alpha \leq N}$, with each ϕ^α providing a smooth approximation of the characteristic function of the given phase, i.e. with $\phi^\alpha = 0$ indicating that one is outside

the phase and $\phi^\alpha = 1$ indicating that one is within the phase². The requirement for a smooth transition between these two cases enforces the existence of a transition region or **interface region** whenever two or more phases meet, i.e. mixed regions where several phases are considered to coexist. Within the coexistence regions, it is quite natural to interpret ϕ^α as representing the individual phase or volume fractions, and thus to impose the sum condition $\sum_{\alpha=1}^N \phi^\alpha = 1$. In addition, it may be desirable or even necessary to, as in the two-phase case, impose the additional restriction $0 \leq \phi^\alpha \leq 1$ on the phase fractions, i.e. to restrict the phasefield values to lie within the N -dimensional Gibbs-simplex

$$\mathcal{GS}^N = \Sigma_1^N \cap [0, 1]^N = \{\phi \in \mathbb{R}^N : \sum_{\alpha=1}^N \phi^\alpha = 1, 0 \leq \phi^\alpha \leq 1 \forall \alpha\}, \quad (6.1)$$

where Σ_1^N denotes the hyperplane of all N -dimensional vectors whose entries sum up to 1.

Remark 31. Note that the conditions in (6.1) are partially redundant as $0 \leq \phi^\alpha \forall \alpha$ and the sum-constraint automatically imply that $\phi^\alpha = 1 - \sum_{\beta \neq \alpha} \phi^\beta \leq 1 \forall \alpha$. As this redundancy can sometimes lead to unnecessary complications - in particular when investigating Lagrange-multipliers associated with the constraints - it is therefore often convenient to use the equivalent but simpler description of \mathcal{GS}^N through³

$$\mathcal{GS}^N = \Sigma_1^N \cap \mathbb{R}_+^N = \{\phi \in \mathbb{R}^N : \sum_{\alpha=1}^N \phi^\alpha = 1, 0 \leq \phi^\alpha \forall \alpha\}. \quad (6.2)$$

◇

While there are often additional difficulties in the modeling of the physics within these coexistence regions, they are also the basis for the approximation of the surface energies (or, in other settings e.g. the entropy) through a **volume** integral

$$\int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) \, d\mathbf{x} \quad (6.3)$$

over the entire (fixed) domain. This is in contrast to the “sharp-interface” representation where this modeling is not necessary, but the surface energy needs to be evaluated in terms of surface integrals over surfaces whose position in time is usually an unknown to be determined as a part of the problem to be solved. In Equation (6.3), a and w are generalizations of the gradient energy density and bulk potential to the multiphase case, whose role within a variational framework is the same as in the simpler two-phase case, namely a penalizes gradients (indicated by the additional $\nabla \phi$ -argument) while w penalizes values deviating from the bulk-values 0 and 1, and ϵ is a parameter for controlling their relative strenghts and thus the resulting interface width. In order for the integral in Equation (6.3) to provide the desired smooth approximation of the surface energy (or entropy) in a multiphase problem, a and w still need to be chosen in an appropriate manner. As, even for multiphase problems, the surface energies in the sharp interface setting are associated with the interfaces between just two coexisting neighboring regions (with regions where more than two phases meet making up lower-dimensional objects such as lines or points), it is quite natural that the multiphase generalizations are usually based upon various modifications of a simple summation procedure over two-phase interactions, such that, in an idealized situation consisting of flat two-phase interphases alone, one would recover the known results from that simpler setting. This approach still leaves a large degree of flexibility

²Note that in the two-phase setting, where it suffices to use a single order paramter, other choices such as ± 1 or $\pm \frac{1}{2}$ for indicating one or the other phase are also popular. While one could in principle also use other values than 0 and 1 in the multiphase case, this choice turns out to be a particularly convenient one.

³Here it is important to maintain the lower bound 0 instead of the upper one as soon as $N \geq 3$. For example, the vector $\phi = (-0.2, 0.6, 0.6)^T$ would satisfy both the sum-constraint and $\phi^\alpha \leq 1 \forall \alpha$.

in constructing particular functionals, which primarily differ in their treatment of triple- and multiphase regions. Combined with the additional degree of freedom in the precise form of the postulated phasefield dynamics - even if based on a gradient type flow - this has lead to a variety of different phasefield models, each with their own advantages and disadvantages.

The one primarily used in this thesis - and the current default within the **Pace3D**-framework - is based on the formulation originally introduced in [52] and [71] within an entropy framework for solidification problems. Considering for simplicity the corresponding energy formulation, the gradient energy density in the isotropic case is defined by

$$a(\phi, \nabla\phi) = \sum_{\beta>\alpha} \gamma^{\alpha\beta} |\mathbf{q}^{\alpha\beta}(\phi, \nabla\phi)|^2, \quad (6.4)$$

with the **generalized antisymmetric gradient vectors**

$$\mathbf{q}^{\alpha\beta}(\phi, \nabla\phi) = \phi^\alpha \nabla\phi^\beta - \phi^\beta \nabla\phi^\alpha. \quad (6.5)$$

This amounts to a summation over energy contributions stemming from two-phase interactions in terms of the respective $\gamma^{\alpha\beta}$ and $\mathbf{q}^{\alpha\beta}$ corresponding to each α - β -pairing.

Extensions to anisotropic gradient energy densities can then be obtained in much the same way as in the two-phase case by setting

$$a(\phi, \nabla\phi) = \sum_{\beta>\alpha} A^{\alpha\beta}(\mathbf{q}^{\alpha\beta}) = \sum_{\beta>\alpha} \gamma^{\alpha\beta} (a^{\alpha\beta})^2 (\mathbf{q}^{\alpha\beta}) |\mathbf{q}^{\alpha\beta}(\phi, \nabla\phi)|^2 \quad (6.6)$$

and a prefactor $a^{\alpha\beta}$ homogeneous of degree 0 and thus depending only on the orientation of $\mathbf{q}^{\alpha\beta}$, but not its norm (i.e. is a function of the normal vector $\mathbf{n}^{\alpha\beta} = \frac{\mathbf{q}^{\alpha\beta}}{|\mathbf{q}^{\alpha\beta}|}$ only).

Remark 32. Note that, while sometimes convenient as it makes the desired quadratic dependence of a on the norm of the gradient vectors easily visible, the definition of the prefactors as a function of the direction only leads to some tricky issues when the norm of $\mathbf{q}^{\alpha\beta}$ approaches zero as this leaves the normal direction - and thus the $a^{\alpha\beta}$ - undefined unless $a^{\alpha\beta}$ does not actually depend on the direction, i.e. unless the gradient energy density is isotropic.

This does not entail any issues in terms of the “total” a -function or the $A^{\alpha\beta}$, since both a and $\frac{\partial a}{\partial \mathbf{q}^{\alpha\beta}}$ (resp. $A^{\alpha\beta}$ and $\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}$) will converge to zero with $|\mathbf{q}^{\alpha\beta}| \rightarrow 0$ (being homogeneous of degree 2 resp. 1 in $\mathbf{q}^{\alpha\beta}$) regardless of the direction along which one approaches this point, but poses a problem when trying to work with the $a^{\alpha\beta}$ themselves as in [75]. \diamond

As outlined above, a generalized w -formulation could be obtained by introducing an additional summation into the two-phase formulation, i.e. by setting

$$\tilde{w}_{mw}(\phi) = 9 \sum_{\beta>\alpha} \gamma^{\alpha\beta} (\phi^\alpha)^2 (\phi^\beta)^2$$

in the well-case resp.

$$\tilde{w}_{mo}(\phi) = \frac{16}{\pi^2} \sum_{\beta>\alpha} \gamma^{\alpha\beta} \phi^\alpha \phi^\beta$$

in the obstacle-case. Unfortunately, using this simple extension, there is a tendency towards the appearance of additional phases throughout the entire interface region, and in particular also in regions that one would naturally want to consist of just the two phases corresponding to neighboring bulk regions. For this reason, an additional penalty term is used in [52] to associate an additional cost with the presence of more than two phases by setting

$$w_{mw}(\phi) = 9 \sum_{\alpha<\beta} \gamma^{\alpha\beta} (\phi^\alpha)^2 (\phi^\beta)^2 + \sum_{\alpha<\beta<\delta} \gamma^{\alpha\beta\delta} (\phi^\alpha)^2 (\phi^\beta)^2 (\phi^\delta)^2 \quad (6.7)$$

resp.

$$w_{mo}(\phi) = \frac{16}{\pi^2} \sum_{\alpha < \beta} \gamma^{\alpha\beta} \phi^\alpha \phi^\beta + \sum_{\alpha < \beta < \delta} \gamma^{\alpha\beta\delta} \phi^\alpha \phi^\beta \phi^\delta. \quad (6.8)$$

The penalty parameters $\gamma^{\alpha\beta\delta}$ are thus used to add an additional energy contribution over triple-phase interactions with an analogous form as the previous two-phase energy contributions.

Remark 33. Note that within a pure two-phase α - β region with $\phi^\alpha = 1 - \phi^\beta$, $\mathbf{q}^{\alpha\beta}$ reduces to $-\nabla\phi^\alpha = \nabla\phi^\beta$, i.e. the gradient energy distribution in such regions reduces precisely to the same form as in the basic two-phase case. Similarly, as all triple-phase terms vanish within two-phase regions and the summation over the two-phase pairings is based on the expressions obtained by replacing ϕ and $1 - \phi$ with ϕ^α resp. ϕ^β (or vice versa) in the two-phase case, w will also reduce to the desired form. \diamond

Remark 34. There are of course many conceivable ways of generalizing the gradient energy density $a(\nabla\phi)$ from the two-phase setting to the multiphase case. The construction used above as a summation over two-phase interactions can be motivated by the observation that the surface energy densities $\gamma^{\alpha\beta}$ correspond to energies between pairings of phases. In addition, in the anisotropic case, the surface energy associated with an interface only depends on the orientation between the adjacent phases and not the absolute orientation of either one of them. This is closely related to another somewhat delicate issue within multiphase regions, namely that of defining appropriate “normal vectors” and angles between phases (the choice $\mathbf{n} = \pm \frac{\nabla\phi}{|\nabla\phi|}$ being a fairly obvious one in the two-phase case).

While the choice $\mathbf{n}^\alpha = \pm \frac{\nabla\phi^\alpha}{|\nabla\phi^\alpha|}$ is still a natural one for the normal vector associated with the phase α “as a whole”, the a priori intuitive choice of defining the orientation between two phases based directly on two such normals may not necessarily lead to the best results. Considering the simple sharp interface example of three straight interfaces meeting in a triple point, it is clear that the angle between every pair of phases is well-defined and constant everywhere except at a single singular point, at which all angles admit an obvious limit though. Within the diffuse interface setting, this triple point is represented by an entire triple point region with small but finite extension. It is clearly desirable with respect to this example for the definition of the relative orientation between any two phases to remain as stable as possible upon the transition from a two-phase to a three-phase (or, more generally, multiphase) region as this may otherwise lead to non-negligible deviations in the energetics unless the ratio of the interface width to the domain size is chosen sufficiently small⁴. As the vectors $-\frac{\nabla\phi^\alpha}{|\nabla\phi^\alpha|}$ of the individual phases will generally follow a smooth transition from one (approximately) constant normal vector to the other during the transition through the triple-point region, due to an increasing influence of the respective third phase, this may be difficult to achieve based on these vectors themselves.

The choice of the generalized gradient vectors $\mathbf{q}^{\alpha\beta}$ replacing the original $\nabla\phi$ corresponds to one possible way of defining a better suited normal $\mathbf{n}^{\alpha\beta} = \frac{\mathbf{q}^{\alpha\beta}}{\|\mathbf{q}^{\alpha\beta}\|}$ between the individual phase pairings, another popular choice being⁵ $\mathbf{n}^{\alpha\beta} = \frac{\nabla\phi^\beta - \nabla\phi^\alpha}{\|\nabla\phi^\beta - \nabla\phi^\alpha\|}$ (see e.g. [61] and [74]). \diamond

Assuming in line with the purpose of this chapter that any additional energetic contribution are prescribed in the form of a simple weighted interpolation of **given** phase-specific contributions⁶ f^α , another question to be addressed is how the nonlinear interpolation functions⁷ $h(\phi)$

⁴Based on geometric considerations, one would expect the influence of these multiphase regions to decrease very quickly with ϵ . Whereas interfaces translate to regions with a volume $\sim \epsilon$, the volumetric contribution of “edges” and “points” are expected to be roughly of the order $\sim \epsilon^2$ resp. $\sim \epsilon^3$. Nevertheless, these regions can also be associated with high energy densities such as e.g. when modeling an elastic problem with corners and edges where one expects singularities in the solution, making these estimates too optimistic.

⁵Note that the choice of sign here in the second case is such that it agrees with the one for $\mathbf{q}^{\alpha\beta}$ in two-phase regions.

⁶Regarding the motivation of using f^α here see remark 36.

⁷For the simplest linear case $h(\phi) = \phi$ this is obvious.

from the two-phase setting can be generalized while maintaining their principle properties. Denoting the interpolation weights for the individual phases by $h^\alpha(\phi)$, the first such property is that $h^\alpha(\mathbf{e}^\beta) = \delta^{\alpha,\beta}$, i.e. the interpolation weights in a bulk-region should be one for the bulk phase and zero for all other phases. A second desirable property is the monotonic increase of $h^\alpha(\phi)$ with ϕ^α . Finally, the interpolation weights should also satisfy the consistency condition $\sum_\alpha h^\alpha(\phi) = 1$, which is necessary to ensure that the interpolation of any property e which does not in fact depend on the phases, $f^\alpha = f \forall \alpha$, satisfies $\sum_\alpha f^\alpha h^\alpha(\phi) = f$. A simple approach ensuring these properties under very basic conditions is the following (see e.g. also [63]). It consists in taking any monotonically increasing (scalar) function⁸ $\tilde{h}(\phi)$ satisfying $\tilde{h}(0) = 0$ and $\sum_\beta \tilde{h}(\phi^\beta) > 0$ for $\sum_\beta \phi^\beta = 1$ and defining the interpolation weights $h^\alpha(\phi)$ through a normalization by

$$h^\alpha(\phi) = \frac{\tilde{h}(\phi^\alpha)}{\sum_\beta \tilde{h}(\phi^\beta)}. \quad (6.9)$$

Remark 35. The third condition $\sum_\alpha h^\alpha(\phi) = 1$ is then satisfied by construction. Together with $\tilde{h}(0) = 0$, this normalization automatically also guarantees that the bulk weights $h^\alpha(\mathbf{e}^\alpha)$ satisfy $h^\alpha(\mathbf{e}^\alpha) = 1$. That the second condition also holds for all ϕ in the Gibbs-simplex follows from a simple direct calculation showing that

$$\frac{\partial h^\alpha(\phi)}{\partial \phi^\alpha} = \frac{\left(\sum_\beta \tilde{h}(\phi^\beta)\right) \tilde{h}'(\phi^\alpha) - \tilde{h}(\phi^\alpha) \frac{\partial}{\partial \phi^\alpha} \sum_\beta \tilde{h}(\phi^\beta)}{\left(\sum_\beta \tilde{h}(\phi^\beta)\right)^2} = \frac{\sum_{\beta \neq \alpha} \tilde{h}(\phi^\beta)}{\left(\sum_\beta \tilde{h}(\phi^\beta)\right)^2} \tilde{h}'(\phi^\alpha).$$

The denominator is obviously non-negative and $\tilde{h}'(\phi^\alpha) \geq 0$ by monotonicity. Finally, the numerator is also non-negative for $0 \leq \phi^\beta \leq 1$ for all $\beta \neq \alpha$ due to the monotonicity condition together with $\tilde{h}(0) = 0$, therefore showing that $\frac{\partial h^\alpha(\phi)}{\partial \phi^\alpha} \geq 0$. \diamond

Combining the different components above, one obtains a phasefield functional of the form

$$\mathcal{F}_\epsilon(\phi) = \int_\Omega \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + f(\phi) \, d\mathbf{x}, \quad (6.10)$$

which is to be minimized under the relevant constraints on the phasefield ϕ . Besides the (relatively simple) local sum-constraint $\sum_{\alpha=1}^N \phi^\alpha = 1$ in the case of the multi-well potential or its somewhat more complex counterpart in terms of the Gibbs-simplex \mathcal{GS}^N in the multi-obstacle case, a further quite common restriction (considered e.g. in [29]) is to fix the total volume $\int_\Omega \phi^\alpha \, d\mathbf{x} \stackrel{!}{=} V^\alpha$ of one or several phases α .

Remark 36. Here, the choice of the letter f (resp. \mathcal{F}) should not be assigned any deeper physical or thermodynamic meaning. Since the functional in this particularly simple setting has no dependence on any of the standard thermodynamic quantities, this is just a matter of notational convenience here since this leads to a set of equations which is closer to the one in [52] (where the - then actual free energy - enters the phasefield equation through the term $-\frac{1}{T} \frac{\partial f}{\partial \phi}$) and the discussion in section 7.1 which also relies on the actual free energy. \diamond

⁸Some popular examples in the two-phase setting will be listed in section 6.2.

6.2 Two-Phase Problems

Before returning to the multi-phasefield model introduced in the previous section in some additional detail, this section will first discuss a number of pure phasefield applications in a pure two-phase setting. Interpreting the variables ϕ^1 and ϕ^2 as the phase fractions of the first and second phase, it is natural to impose the constraint $\phi^1 + \phi^2 = 1$ as a special case of the one in Section 6.1. This implies that each phase fraction is known as soon as the other one is, such that these problems can, as in the discussion in the introduction, be completely described in terms of a single order parameter ϕ alone, taken here for definiteness as $\phi := \phi^1$.

It is clear that this leads to a number of significant simplifications as compared to the more general multi-phasefield model from the outlined above. Firstly, this reduces the vectorial phasefield equations to a single scalar one for ϕ . Secondly - and this is partially related to the particular phasefield model used here - the generalized gradient vector \mathbf{q}^{12} reduces back to the (negative of) the basic gradient vector $-\nabla\phi$, i.e. unless the gradient energy density is anisotropic, there is no nonlinearity in the “spatial” part of the functional. Thirdly, unless additional constraints are present, the representation of the admissible set becomes essentially trivial. As the sum-constraint is already fully incorporated into the formulation, it simply drops out in the double-well case, while the restriction to the Gibbs-simplex in the double-obstacle case reduces to the box-constraint $\phi(\mathbf{x}) \in [0, 1]$.

Remark 37. Even though it may seem that the interest of the two-phase setting (as a particularly simple special case of the more general multiphase-setting) is quite limited, this is somewhat misleading. On the one hand, there are in fact a number of interesting applications of the phasefield method (such as capillary rise problems and topology optimization problems) which, while not a priori restricted to two phases, in practice often arise naturally in precisely this setting. On the other hand, even though the simplifications above as compared to the multiphase models are certainly significant, many of the essential challenges arising in the general case (high resolution requirements due to the small length-scale ϵ , lack of convexity, bound-constraints, treatment of strong anisotropies, ...) are actually inherited - though possibly in an exacerbated fashion - from the ones in the two-phase case⁹. \diamond

The problems to be considered in this section will, similar to Section 3, be based on slight variations of the functional

$$\mathcal{F}_\epsilon(\phi) = \int_{\Omega} \epsilon a(\nabla\phi) + \frac{1}{\epsilon} w(\phi) + f(\phi) \, d\mathbf{x}, \quad (6.11)$$

or, with a in the simplest case given by $a(\nabla\phi) = \gamma|\nabla\phi|^2$, its isotropic version

$$\mathcal{F}_\epsilon(\phi) = \int_{\Omega} \gamma\epsilon|\nabla\phi|^2 + \frac{1}{\epsilon} w(\phi) + f(\phi) \, d\mathbf{x}. \quad (6.12)$$

The bulk potential may be either given by the **double well potential**

$$w_{dw}(\phi) = 9\gamma\phi^2(1 - \phi)^2 \quad (6.13)$$

or the **double obstacle potential**

$$\tilde{w}_{ob}(\phi) = w_{ob}(\phi) + 1_{[0,1]}(\phi), \quad (6.14)$$

where

$$w_{ob}(\phi) = \frac{16}{\pi^2} \gamma \phi(1 - \phi), \quad (6.15)$$

⁹A notable exception here is given by the triple-phase terms in the bulk potential w , which, depending on the size of the penalty parameters, may lead to an additional stiffness in the equations absent in the presence of only two phases.

is the smooth part of the obstacle potential and $1_{[0,1]}$ is the indicator function

$$1_{[0,1]}(\phi) = \begin{cases} 0 & , 0 \leq \phi \leq 1, \\ +\infty & , \text{else.} \end{cases}$$

of the interval $[0, 1]$, whose role it is to maintain the phasefield values within the interval $[0, 1]$. Alternatively, one can also impose this restriction directly in terms of an additional constraint $0 \leq \phi \leq 1$, and, since $1_{[0,1]}$ vanishes for all such values, work with the smooth part $w_{ob}(\phi)$ only instead of $\tilde{w}_{ob}(\phi)$.

Remark 38. Whereas the restriction of the ϕ -values to lie within $[0, 1]$ is a necessity when using w_{ob} , in the case of the double-well potential one may or may not choose to impose this additional constraint. As these bounds entail some notable complications in the solution of the resulting phasefield problems below, which would only seem to be counterbalanced by large “pure” bulk regions (i.e. regions with $\phi = 0$ or $\phi = 1$) when using double obstacle potential, this possibility will not further be considered here¹⁰. \diamond

In the simplest cases, $f(\phi)$ is given by an interpolation

$$f(\phi) = \Delta f h(\phi) = (f^1 - f^2)h(\phi), \quad (6.16)$$

describing the energy difference $\Delta f = f^1 - f^2$ associated with transforming from the first to the second phase. $h(\phi)$ is some appropriate interpolation function, whose choice provides some additional freedom in modeling the energy contributions within two-phase regions.

Classical choices¹¹ (see e.g. [52]) together with their derivatives are shown in Table 6.2. Unlike h_0 , for which the contribution of the f^α for each phase are directly proportional to the phase fractions, the energy contributions for the other two functions remain closer to that of the predominant phase, with a steeper transition around $\phi = \frac{1}{2}$ and therefore a sharper transition region in an energetic sense than the phasefield variable itself.

	$h(\phi)$	$h'(\phi)$	$h''(\phi)$
h_0	ϕ	1	0
h_1	$\phi^2(3 - 2\phi)$	$6\phi(1 - \phi)$	$6 - 12\phi$
h_2	$h_2(\phi) = \phi^3(6\phi^2 - 15\phi + 10)$	$30\phi^2(1 - \phi)^2$	$60\phi(1 - \phi)(1 - 2\phi)$

6.2.1 The Steady-State and Dynamic Phasefield Equations

From the functional \mathcal{F}_ϵ above and using the appropriate version of the functional derivative $\frac{\delta \mathcal{F}_\epsilon}{\delta \phi} = \frac{\partial \mathcal{F}}{\partial \phi} - \nabla \cdot \left(\frac{\partial \mathcal{F}_\epsilon}{\partial \nabla \phi} \right)$, it is now seemingly straightforward to (re)derive the first-order necessary conditions to be satisfied by any potential minimizer ϕ of \mathcal{F}_ϵ and, based upon the postulate of a gradient flow and under some additional assumptions, their “dynamic” version.

At least in the obstacle case, there is an additional constraint in terms of $\phi \in [0, 1]$ though. From this and the discussion in Section 4.3 one expects the appearance of two additional multipliers in the resulting phasefield equation. Even though these do have a very fundamental impact, they are often simply left out in the more applied literature.

¹⁰For coupled problems though, allowing negative ϕ -values in combination with simple interpolation schemes for the material properties (e.g. weighing the material properties by the ϕ -values of the respective phase) may lead to very serious issues such as e.g. negative diffusivities or stiffnesses. One way of dealing with this is of course to simply add the bound-constraints on the phasefield variable. Another one - avoiding the additional complications this constraint entails - is to modify the basic interpolation schemes outside the range $[0, 1]$ such as to ensure physically valid values of the interpolated quantities.

¹¹In relation with the variational approach, it should be noted that $H^1(\Omega) \hookrightarrow L^q$ with $1 \leq q < \infty$ if $n = 2$ and $1 \leq q \leq 6$ if $n = 3$, i.e. integrability is ensured without knowing a priori whether ϕ is particularly well-behaved or not.

Remark 39. This may partly be due to the fact that their action can in principle be “hidden” behind an appropriate projection operator¹², but the practical usefulness of this approach is highly problem-dependent. While this can be very convenient, in particular when using an explicit time-stepping scheme involving a purely local projection operations, this need not be the case anymore when using other time-discretizations or if one adds additional constraints (see e.g. Remark 42). ◇

Since analyzing the impact of box-constraint is - at least at a purely formal level disregarding regularity issues and the precise meaning of this constraint¹³ - mostly basic, it therefore seems worthwhile to at least indicate the basic reasoning from which the resulting steady-state equations and their complete form can in principle be derived. To keep the discussion as simple as possible, it will in addition be assumed that $a(\nabla\phi) = \gamma|\nabla\phi|^2$, i.e. that the phasefield energy is isotropic¹⁴. Taking the directional derivative of the phasefield functional \mathcal{F}_ϵ leads to the first-order necessary condition

$$\mathcal{F}'_\epsilon(\phi; \psi) = \int_{\Omega} 2\gamma\epsilon\nabla\phi \cdot \nabla\psi + \frac{1}{\epsilon}w'(\phi)\psi + f'(\phi)\psi \, d\mathbf{x} \geq 0,$$

or, assuming sufficient regularity, an immediate integration by parts modifies this into

$$\mathcal{F}'_\epsilon(\phi; \psi) = \int_{\partial\Omega} 2\gamma\epsilon\frac{\partial\phi}{\partial\mathbf{n}}\psi \, ds + \int_{\Omega} \left(-2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) \right)\psi \, d\mathbf{x} \geq 0$$

for all admissible directions ψ . At this point, there is now a strong difference between the (unconstrained when expressed in the reduced form) double-well case and the double-obstacle one. In the former case, ψ is (up to some regularity requirements due to the underlying function space) essentially arbitrary. Since in particular both ψ and $-\psi$ are admissible as long as ψ is so, one has both $\mathcal{F}_\epsilon(\phi; \psi) \geq 0$ and $\mathcal{F}_\epsilon(\phi; -\psi) = -\mathcal{F}_\epsilon(\phi; \psi) \geq 0$, and thus

$$\mathcal{F}'_\epsilon(\phi; \psi) = \int_{\partial\Omega} 2\gamma\epsilon\frac{\partial\phi}{\partial\mathbf{n}}\psi \, ds + \int_{\Omega} \left(-2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) \right)\psi \, d\mathbf{x} = 0 \quad \forall \psi \in H^1(\Omega). \quad (6.17)$$

From this it is, at least formally, quite easy to obtain a more explicit characterization of the minimizer. First “testing” this equality with functions ψ vanishing on the boundary, one can conclude that one must have $-2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) = 0$ in Ω for this to hold. Provided this is the case, the volume integral vanishes in Equation (6.17), and by varying the boundary values of ψ , it then follows that $\frac{\partial\phi}{\partial\mathbf{n}}$ will also have to be zero. In summary, one thus has to satisfy

$$\begin{cases} -2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) = 0 & \text{in } \Omega, \\ \frac{\partial\phi}{\partial\mathbf{n}} = 0 & \text{on } \partial\Omega, \end{cases} \quad (6.18)$$

or, in the general anisotropic case,

$$\begin{cases} -\epsilon\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi} \right) + \frac{1}{\epsilon}w'(\phi) + f'(\phi) = 0 & \text{in } \Omega, \\ \frac{\partial a}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega. \end{cases} \quad (6.19)$$

In contrast, when subject to the box-constraints, ψ and $-\psi$ are usually not both admissible directions¹⁵, and even a formal argument requires a little more care.

¹²This topic will be discussed in some more detail in Section 6.3, where, even in the discrete and fully explicit case, this projection can become somewhat more complex depending on the precise form of dynamics chosen.

¹³On a less formal level, dealing with an obstacle-type constraint is in fact a quite involved question. This will not be discussed in any detail here, but some of the more relevant issues will nevertheless be pointed out in Remark 40 below.

¹⁴This is mostly irrelevant, but avoids a tedious discussion for relating the signs of $\frac{\partial a}{\partial \mathbf{n}}$ and $\frac{\partial \phi}{\partial \mathbf{n}}$.

¹⁵Unlike for the double-well case without constraints or e.g. equality constraints such as the volume-constraints, the admissible directions in combination with the inequality constraints do not form a linear space but a cone (see Chapter 4.3 for a basic discussion).

If ϕ is at least locally continuous and $0 < \phi < 1$ at some point inside the domain, there would, by this continuity, exist some neighborhood on which ϕ would continue to satisfy this inequality and one could thus take (a sufficient amount of) local variations ψ for which both ψ and $-\psi$ would be admissible. From this one can deduce that

$$-2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) = 0 \text{ if } 0 < \phi < 1.$$

If one instead has $\phi = 0$ in some region, the only admissible variations ψ will locally have to satisfy (in the appropriate sense) $\psi \geq 0$, from which it then only follows that

$$\mu^- := -2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) \geq 0.$$

Applying the same argument to those “regions” where $\phi = 1$, it similarly follows that

$$-\mu^+ := -2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) \leq 0.$$

As both μ^+ and μ^- are non-negative, these three conditions can be summarized to

$$-2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) + \mu^+ - \mu^- = 0 \tag{6.20}$$

combined with the complementarity conditions

$$\mu^+(1 - \phi) = 0 \quad \text{and} \quad \mu^-\phi = 0 \tag{6.21}$$

enforcing that μ^+ resp. μ^- has to vanish unless the associated constraint is active, i.e. $1 - \phi = 0$ or $\phi = 0$.

That this conclusion is not quite as simple as in the unconstrained case also leads to some complications when trying to “focus” on the boundary. In fact, in contrast to the previous argument, it does not follow that the volume integral simply drops out. Instead, using Equation (6.20), one can only conclude that

$$\mathcal{F}'_\epsilon(\phi; \psi) = \int_{\partial\Omega} 2\gamma\epsilon \frac{\partial\phi}{\partial\mathbf{n}} \psi \, ds + \int_{\Omega} (\mu^- - \mu^+) \psi \, dx \geq 0$$

with $\mu^\pm \geq 0$ together with the additional restrictions in Equation (6.21). This in principle also makes the argument pertaining to the correct boundary conditions a little more complicated. The simplest case is again the one if there is some neighborhood of the boundary where $0 < \phi < 1$, as then both μ^+ and μ^- vanish in this region and one can additionally (locally) choose ψ freely, leading, as before, to $\frac{\partial\phi}{\partial\mathbf{n}} = 0$. If in contrast $\phi = 0$ or $\phi = 1$ in some neighborhood of the boundary, either μ^+ or μ^- need not (and in combination with w_{ob} normally will not) vanish. Nevertheless, since this means that one locally has a bulk-region, there is no need to pass through an indirect argument using the test-functions as this implies through an obvious limit that $\nabla\phi = 0$ on the boundary and thus obviously also $\frac{\partial\phi}{\partial\mathbf{n}} = 0$.

What is tricky though is the case when the boundary is not locally in a bulk-region, but ϕ tends to one of the critical values 0 or 1 as one approaches the boundary. Even though μ^+ and μ^- and therefore the volume-integral would locally vanish, one would face a sign restriction on the variations of ψ on the boundary¹⁶ and would thus a priori only be able to conclude that $\frac{\partial\phi}{\partial\mathbf{n}} \geq 0$

¹⁶Note that this does not follow from e.g. the constraint $0 \leq \phi \leq 1$ a.e. in Ω by itself, since the boundary is a set of measure zero. It is only in combination with the additional regularity - such as the one imposed through the space $H^1(\Omega)$ - that this as well as talking about the values of ψ on the boundary becomes meaningful.

if $\phi \rightarrow 0$ resp. $\frac{\partial \phi}{\partial \mathbf{n}} \leq 0$ if $\phi \rightarrow 1$ as one approaches the boundary¹⁷. A strict inequality in these conclusions is in contradiction with the assumptions on the profile though, as this to first order means that any point in the boundary with $\phi = 0$ would have to be approached from **below** and any point with $\phi = 1$ from **above**. Combining all these observations, the first-order necessary conditions can finally be summarized to

$$\begin{cases} -2\gamma\epsilon\Delta\phi + \frac{1}{\epsilon}w'(\phi) + f'(\phi) + \mu^+ - \mu^- = 0 & \text{in } \Omega, \\ 0 \leq \phi \leq 1, \mu^\pm \geq 0, \mu^+(1-\phi) = 0 \text{ and } \mu^-\phi = 0, & \\ \frac{\partial \phi}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega, \end{cases} \quad (6.22)$$

respectively in the general potentially anisotropic case

$$\begin{cases} -\epsilon\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi} \right) + \frac{1}{\epsilon}w'(\phi) + f'(\phi) + \mu^+ - \mu^- = 0 & \text{in } \Omega, \\ 0 \leq \phi \leq 1, \mu^\pm \geq 0, \mu^+(1-\phi) = 0 \text{ and } \mu^-\phi = 0. & \\ \frac{\partial a}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega, \end{cases} \quad (6.23)$$

Remark 40. It needs to be stressed that, while in principle consistent, the argument above is purely formal. Making it rigorous requires quite a bit of additional work and raises a number of - often also practically important - questions.

Besides the derivation itself, one central question is of course the sense in which the various quantities and equations should be interpreted and is thus essentially a matter of function spaces. This point is more delicate in problems involving obstacles than for other more “global” constraints or unconstrained problems. In particular, whereas the solutions of many classical equations in the applied physics based on a Laplace-type operator are known to be as regular as the data permits (and in particular C^∞ in the interior of the domain if the right-hand side is so), the same is not true for obstacle-type problems. In fact, even though both obstacles here are given by constants (and are thus extremely regular), even the “optimal” unperturbed one-dimensional phasefield profile has, as will be recalled in Remark 47, a jump in its second derivative upon crossing from the bulk to an interface region. For this reason, it is clear that one should certainly not expect for Equation (6.23) to hold in the classical sense.

An also in practice very important question is therefore whether it does at least hold in the strong (a.e.) sense in e.g. $L^2(\Omega)$, or whether it has to be interpreted in an even weaker sense. Even though $\phi \in H^1(\Omega)$ does at least (the trace-operator being well-defined) ensure a “continuity through hypersurfaces”, the same does not a priori hold for the gradient of ϕ . On the one hand, jumps in the gradient of ϕ would be expected to significantly increase numerical errors due to a reduction in the convergence order of any standard numerical scheme. On the other hand, Equation (6.23) could then, through the second derivative in its first term, contain surface measures. As both w' and any reasonably chosen f' are bounded for $0 \leq \phi \leq 1$, this contribution could only be counterbalanced by equally measure-valued multipliers μ^\pm . While this has no precise meaning in the discrete case, any convergent numerical scheme relying (explicitly or implicitly) on these multipliers would have to mimic this behavior as the discretization is refined. This requires for the discrete counterpart of μ^\pm to locally tend to ∞ in some sense, and is thus highly likely to cause additional numerical difficulties.

Fortunately, both theoretical predictions (see e.g. [28] for some regularity results for obstacle problems) and computational practice¹⁸ indicate that this issue does usually not arise. Nevertheless, the treatment of the transition region between the interface and the bulk remains a

¹⁷This is a technically somewhat difficult point as $\frac{\partial \phi}{\partial \mathbf{n}}$ need not be “intrinsically” defined on hypersurfaces as $\phi \in H^1(\Omega)$ by itself does not provide sufficient regularity for this. A simple example illustration of this in the one-dimensional case will be discussed in Remark 46 (the interested reader is referred to [48] for a much more general discussion of this topic). The central assumption made here is that μ^\pm are both integrable functions and not measure-valued.

¹⁸In the sense that calculated phasefield profiles do indeed exhibit a behavior consistent with a continuous transition of $\nabla \phi$ to $\mathbf{0}$ as one approaches the bulk.

somewhat delicate point from a numerical point of view¹⁹.

For some related discussion in the phasefield context, the reader is referred to e.g. [60]. More general background within a more abstract setting (including e.g. also the question of the regularity of free boundaries) can, among many other sources, be found in [43], [28], [37] and [6] and the references therein. \diamond

Remark 41. It should also be noted that, within the bulk-regions, if $f'(0) = f'(1) = 0$, μ^+ resp. μ^- are given by the constant value $\frac{16\gamma}{\pi^2\epsilon}$ as the only “active” term in Equation (6.23) is given by the derivative of the bulk-potential w in 0 or 1²⁰. In order to eliminate this dependence on ϵ , it can therefore sometimes be convenient (as in [60]) to simply rescale the multipliers by instead using $\hat{\mu}^\pm := \epsilon\mu^\pm$, i.e. modifying Equation (6.23) to

$$-\epsilon\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi} \right) + \frac{1}{\epsilon} w'(\phi) + f'(\phi) + \frac{1}{\epsilon} (\hat{\mu}^+ - \hat{\mu}^-) = 0$$

subject to the same complementarity conditions as μ^\pm . \diamond

Remark 42. If additional constraints are added, the situation also becomes much more involved at a practical level, a very common and quite representative example being the volume-constraint on the phases. Even though this is in principle quite simple at a theoretical level - the integral constraint being very robust with respect to e.g. regularity issues and requiring only a very low-dimensional multiplier - the phasefield equation now involves both a nonlinear but **local** inequality constraint, and a linear but **global** constraint on the phasefield values.

This is in particular a difficulty for numerical schemes avoiding the use of multipliers through a projection-based approach. While the projection operations associated with each constraint separately are relatively straightforward to deal with (at least provided the outer algorithm can make use of the strict locality of the box-constraint), the actually required projection would be a both global and nonlinear one, leading to a very complex and expensive operation.

In practical terms, this can be dealt with in a quite satisfactory manner by instead satisfying a slight modification of the actual first-order necessary condition introducing an additional modified driving force as in [53] (see also [29]). Through an additional simplification (in the obstacle case), this can be turned into a quite notably more efficient algorithm based on purely sequential operations. It should be kept in mind though that these modification are primarily algorithmically motivated, and the solution thus corresponds to a (usually very accurate for the time-stepping scheme actually used in [53]) approximation of the one minimizing the underlying functional. \diamond

In contrast to the conditions characterizing a minimizer of \mathcal{F}_ϵ in the “steady-state” case above, the dynamics of the problem are not directly implied by the functional and the constraints alone. The standard approach in the Allen-Cahn case is that of a (potentially weighted) L^2 -gradient flow, i.e. in the unconstrained case one posulates $\frac{\partial \phi}{\partial t} \sim -\frac{\delta \mathcal{F}_\epsilon}{\delta \phi}$. Assuming for the proportionality to be of the form $\tau\epsilon$ and combined with the expression for the gradient in Equation (6.19), one obtains the dynamic counterpart

$$\tau\epsilon \frac{\partial \phi}{\partial t} = \epsilon\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi} \right) - \frac{1}{\epsilon} w'(\phi) - f'(\phi) \quad (6.24)$$

to Equation (6.19) in the double-well case.

Remark 43. The prefactor $\tau\epsilon$ is the “default” version within the **Pace3D**-framework, with τ potentially depending on e.g. $\nabla\phi$ in order to also enable the incorporation of anisotropic effects into the dynamics in addition to the ones in the gradient energy density. The additional use of

¹⁹One consequence of the expected jump in the second derivative being that one does not expect a formally fourth (or even higher order) scheme to actually achieve this increase in accuracy.

²⁰Even if f' does not vanish for $\phi \in \{0, 1\}$ - as when using $h(\phi) = \phi$ - the other term is generally largely dominant.

ϵ in contrast does not by itself have any direct effect in this context as one could equally well rescale time. As an alternative to the use of τ , one could also use its inverse $m := \frac{1}{\tau}$ as a prefactor on the right-hand side. In contrast to the multiphase situation in Section 6.3, this distinction is of little interest for two-phase problems. \diamond

The situation is again somewhat more complicated when using the obstacle-potential (or more precisely, when enforcing the box-constraint on ϕ) as it naturally leads to counter-forces for any change in the phasefield which would lead outside the admissible set. Arguing either in a manner similar to the one for the steady-state case or simply “generalizing” the expression in Equation (6.23) by adding a time-derivative compatible with the FONC, one obtains the evolution equation

$$\begin{cases} \tau\epsilon \frac{\partial\phi}{\partial t} = \epsilon\nabla \cdot \left(\frac{\partial a}{\partial \nabla\phi} \right) - \frac{1}{\epsilon} w'(\phi) - f'(\phi) - (\mu^+ - \mu^-), \\ 0 \leq \phi \leq 1, \mu^\pm \geq 0, \mu^+(1-\phi) = 0 \text{ and } \mu^-\phi = 0. \end{cases} \quad (6.25)$$

Remark 44. Both Equations (6.24) and (6.25) still need to be complemented, as before, with the boundary condition $\frac{\partial a}{\partial n} = 0$. \diamond

6.2.2 The One-Dimensional Case and the “Standard” Phasefield Profiles

Before continuing the discussion above, it is instructive to take a look at the simple (parameterized families of) analytical solutions for the steady-state problems in the well and obstacle case in their simplest one-dimensional form, as these on the one hand indicate some qualitative properties expected in the more general case and on the other hand are the profiles underlying the link between the diffuse and sharp interface models as well as the thin interface asymptotics. In this setting, one is thus looking for a function ϕ satisfying

$$-2\gamma\epsilon \frac{d^2\phi}{dx^2} + \frac{1}{\epsilon} \frac{dw}{d\phi}(\phi) + \frac{df}{d\phi}(\phi) = 0 \quad (6.26)$$

within the interface region²¹, with w being either the double-well or double-obstacle bulk potential. An important difference between these two cases is that $\frac{\partial w_{do}}{\partial\phi} = \frac{16\gamma}{\pi^2}(1-2\phi)$ is an affine function - thus allowing for applying the standard linear theory, at least within the interface region and provided $\frac{\partial f}{\partial\phi}$ is sufficiently simple - whereas $\frac{\partial w_{dw}}{\partial\phi}$ is a cubic function of ϕ , and requires other tools for determining an analytical solution.

The standard approach for the latter case is based on the existence of a first integral for (6.26), a fact which has its own inherent interest from an energetic point of view also in the obstacle case. More precisely, a multiplication of Equation (6.26) with $\frac{d\phi}{dx}$ and using²² $\frac{d^2\phi}{dx^2} \frac{d\phi}{dx} = \frac{d}{dx} \frac{1}{2} \left| \frac{d\phi}{dx} \right|^2$ as well as $\frac{dw}{d\phi}(\phi) \frac{d\phi}{dx} = \frac{d}{dx} w(\phi)$ and $\frac{df}{d\phi}(\phi) \frac{d\phi}{dx} = \frac{d}{dx} f(\phi)$ ²³ shows that

$$\left(-2\gamma\epsilon \frac{d^2\phi}{dx^2} \frac{1}{\epsilon} \frac{dw}{d\phi}(\phi) + \frac{df}{d\phi}(\phi) \right) \frac{d\phi}{dx} = \frac{d}{dx} \left(-\gamma\epsilon \left| \frac{d\phi}{dx} \right|^2 + \frac{1}{\epsilon} w(\phi) + f(\phi) \right) = 0$$

and thus that

$$-\gamma\epsilon \left(\frac{d\phi}{dx} \right)^2 + \frac{1}{\epsilon} w(\phi) + f(\phi) = \text{const} \quad (6.27)$$

²¹The restriction to this region avoids the issue of the multipliers μ^\pm in the obstacle case.

²²More generally, one could use the relation $-\left(\frac{d}{dx} \frac{\partial a}{\partial \left(\frac{d\phi}{dx} \right)} \right) \frac{d\phi}{dx} = \frac{d}{dx} \left(a \left(\frac{d\phi}{dx} \right) - \frac{\partial a}{\partial \left(\frac{d\phi}{dx} \right)} \frac{d\phi}{dx} \right)$. The use of an anisotropic a -function seems of little interest in the one-dimensional case though.

²³Valid as long as the only dependence on x in w resp. f arises solely through that of $\phi(x)$.

within any region where one does not have to consider the bound constraints, i.e. within the region $0 < \phi(x) < 1$ in the obstacle case and everywhere in the domain in the double-well case. Besides the very useful information contained in Equation (6.27) from an energetic point of view, an advantage (at least in the nonlinear double-well case) over the original second-order differential equation is that it directly implies

$$\sqrt{\gamma\epsilon} \frac{d\phi}{dx} = \pm \sqrt{\frac{1}{\epsilon} w(\phi) + f(\phi) - \text{const.}}$$

This gives, within each region of monotonicity of ϕ , rise to a separable first-order ODE and thus turns the problem of determining a solution to the original equation into the potentially simpler one of determining an indefinite integral, whose solution one can then try to match to the appropriate boundary conditions.

For some particular cases, this can in fact be done analytically. Assuming for example that w is given by the double-well potential, that there is no driving force and that the domain is an infinite one, it turns out that the constant in Equation (6.27) can be taken to be zero, and that the solutions to the resulting nonlinear ODE $\sqrt{\gamma\epsilon} \frac{d\phi}{dx} = \sqrt{\frac{1}{\epsilon} w(\phi)}$ are known to be given by²⁴

$$\phi(x) = \frac{1}{2} \left(1 \pm \tanh \left(\frac{3}{2\epsilon} x + c \right) \right) \quad (6.28)$$

valid for $-\infty < x < +\infty$.

Remark 45. In practical computations, one of course never deals with actually infinite domains. Even though the solution (6.28) is then not valid, the deviations from this ideal profile are usually very small due to the exponential convergence to the bulk-values 0 and 1. Similarly, in higher-dimensional and multi-interface and/or multiphase simulations involving several interfaces, these “tails” will in principle lead to an interaction between all phases and all interfaces. As long as there is a reasonable distance between the small but “crucial” parts of the transition regions, the energetic effect of these long-range interactions is typically negligible as compared to other sources of deviations such as curvature-induced changes, discretization errors and the changes in the energetics and profiles induced by multiphase regions. \diamond

The situation in the double-obstacle case is both simpler and more complex. In fact, using the expression for $\frac{dw_{d\phi}}{d\phi}$, the original second-order PDE (6.26) reduces, provided $\frac{df}{d\phi}$ vanishes, to $\frac{d^2\phi}{dx^2} + \frac{16}{\pi^2\epsilon^2}\phi = \frac{8}{\pi^2\epsilon^2}$ within the interface region. It is well-known that the general solution to this is equation is given by the superposition of a particular solution $\phi_p(x)$, which can here be chosen as $\phi_p(x) = \frac{1}{2}$, and a homogeneous solution of the form $\phi_h(x) = \tilde{c}_1 \sin\left(\frac{4}{\pi\epsilon}x\right) + \tilde{c}_2 \cos\left(\frac{4}{\pi\epsilon}x\right)$ corresponding to the two roots of the characteristic equation $\lambda^2 + \frac{16}{\pi^2\epsilon^2} = 0$. Appropriately choosing the constants, this can equivalently be rewritten as

$$\phi(x) = \frac{1}{2} \left(1 \pm \tilde{c}_1 \sin\left(\frac{4}{\pi\epsilon}x + \tilde{c}_2\right) \right). \quad (6.29)$$

²⁴One way of obtaining this is through a partial fraction decomposition. Inserting the definition of w and restricting the focus to the range $0 < \phi < 1$ leads to $\frac{d\phi}{dx} = \pm \frac{3}{\epsilon} \phi(1 - \phi)$. After a separation of the variables, ϕ thus has to satisfy

$$\int \frac{d\phi}{\phi(1-\phi)} = \int \frac{1}{\phi} + \frac{1}{1-\phi} d\phi = \ln(\phi) + \ln(1-\phi) = \ln\left(\frac{\phi}{1-\phi}\right) = \pm \frac{3}{\epsilon} \int dx + \tilde{c} = \pm \frac{3x}{\epsilon} + \tilde{c}.$$

Taking exponentials on both sides simplifies this to $\frac{\phi}{1-\phi} = e^{\pm \frac{3x}{\epsilon} + \tilde{c}}$, or, solving for ϕ ,

$$\phi = \frac{e^{\pm \frac{3x}{\epsilon} + \tilde{c}}}{ce^{\pm \frac{3x}{\epsilon} + \tilde{c}} + 1} = \frac{1}{2} \left(\frac{e^{\pm \frac{3x}{\epsilon} + \tilde{c}} + 1}{ce^{\pm \frac{3x}{\epsilon} + \tilde{c}} + 1} + \frac{e^{\pm \frac{3x}{\epsilon} + \tilde{c}} - 1}{ce^{\pm \frac{3x}{\epsilon} + \tilde{c}} + 1} \right) = \frac{1}{2} \left(1 + \frac{e^{\pm \frac{3x}{\epsilon} + \tilde{c}} - 1}{ce^{\pm \frac{3x}{\epsilon} + \tilde{c}} + 1} \right).$$

From this, the result follows in combination with $\tanh(y) = \frac{e^y - e^{-y}}{e^y + e^{-y}} = \frac{e^{2y} - 1}{e^{2y} + 1}$, defining $c = \pm \frac{\tilde{c}}{2}$ and using the antisymmetry of the tanh-function.

In order to obtain an acceptable solution describing the transition between two bulk-regions, it remains to adjust the constants \tilde{c}_1 and \tilde{c}_2 . The right guess ensuring this (renaming \tilde{c}_2 to c for consistency with the double-well solution in Equation (6.28)) is to set

$$\phi(x) = \frac{1}{2} \left(1 \pm \sin \left(\frac{4}{\pi\epsilon} x + c \right) \right), \quad x \in \left(-c - \frac{\pi^2}{8}\epsilon, -c + \frac{\pi^2}{8}\epsilon \right), \quad (6.30)$$

prolongated by the respective bulk values 0 and 1 outside this inner region.

Remark 46. Even though it is intuitively clear that the profile in Equation (6.30) is indeed the correct one as it is the one matching the “most smoothly” with the bulk-regions, this is not at all obvious based on Equation (6.26) and the condition of having to reach the values of 0 and 1. In fact, the (here implicit) condition that ϕ should at least lie in H^1 implies, in the one-dimensional setting, that ϕ at least has to be continuous. Based on Equation (6.29), this can be achieved using any constant $\tilde{c}_1 \geq 1$ and then simply truncating the profile to the bulk-values at the points where it starts going outside the bounds. Fixing for convenience the constant \tilde{c}_2 to zero (corresponding to centering the interface at $x = 0$) and choosing for definiteness the increasing profile, this still leaves a one-parameter family of potential solutions $\phi(x) = \frac{1}{2} \left(1 + \tilde{c}_1 \sin \left(\frac{4}{\pi\epsilon} x \right) \right)$ based on the ODE withing the interface and the continuity condition alone, and choosing the correct one therefore requires some additional information.

The safest way to argue here is by simply considering the total phasefield energy as function of the parameterization in terms of \tilde{c}_1 and to choose the relevant solution accordingly. A simple calculation shows that, given $\tilde{c}_1 \geq 1$, the interface defined as the region between the points where the profile first touches 0 resp. 1 will then extend from $x_l := -\frac{\epsilon\pi}{4} \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right)$ to $x_r := \frac{\epsilon\pi}{4} \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right)$. Further inserting this profile into the definitions of the gradient and bulk energy densities leads to

$$\epsilon a \left(\frac{d\phi}{dx} \right) = \gamma \frac{4}{\pi^2 \epsilon} \tilde{c}_1^2 \cos^2 \left(\frac{4}{\pi\epsilon} x \right) \quad \text{and} \quad \frac{1}{\epsilon} w(\phi) = \gamma \frac{4}{\pi^2 \epsilon} \left(1 - \tilde{c}_1^2 \sin^2 \left(\frac{4}{\pi\epsilon} x \right) \right),$$

based on which a simple calculation²⁵ shows that

$$\mathcal{F}_\epsilon(\tilde{c}_1) = \int_{x_l}^{x_r} \epsilon a \left(\frac{d\phi}{dx} \right) + \frac{1}{\epsilon} w(\phi) dx = \frac{2\gamma}{\pi} \left(\sqrt{\tilde{c}_1^2 - 1} + \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right) \right).$$

Even though the first term is increasing in \tilde{c}_1 and the second one decreasing, differentiation this expression shows that

$$\frac{d\mathcal{F}_\epsilon(\tilde{c}_1)}{d\tilde{c}_1} = \frac{2\gamma}{\pi} \frac{\tilde{c}_1 - \frac{1}{\tilde{c}_1}}{\sqrt{\tilde{c}_1^2 - 1}},$$

i.e. the contribution by the first term is dominant and the energy is increasing for all $\tilde{c}_1 > 1$. The lowest energy is thus indeed achieved for the minimal admissible $\tilde{c}_1 = 1$, i.e. the previously obtained profile.

This is also a somewhat tedious argument as it in particular relies heavily on explicitly calculating the total energy. A much faster argument is to proceed analogously as in the previous section, which is in the current case much simpler to make rigorous as there are no issues with the

²⁵Adding both contributions and using $\cos^2(\theta) - \sin^2(\theta) = \cos(2\theta)$ shows that the integral can be rewritten as

$$\mathcal{F}_\epsilon(\tilde{c}_1) = \frac{4\gamma}{\pi^2 \epsilon} \int_{x_l}^{x_r} \left(\tilde{c}_1^2 \cos \left(\frac{8}{\pi\epsilon} x \right) + 1 \right) dx = \frac{\gamma}{\pi} \left(\tilde{c}_1^2 \sin \left(2 \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right) \right) + 2 \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right) \right).$$

This simplifies further by inserting $\sin(2\theta) = 2 \sin(\theta) \cos(\theta)$ to get $\tilde{c}_1^2 \sin \left(2 \sin^{-1} \left(\frac{1}{\tilde{c}_1} \right) \right) = 2\tilde{c}_1 \cos \left(\sin^{-1} \left(\frac{1}{\tilde{c}_1} \right) \right)$, from which the result follows in combination with $\cos \left(\sin^{-1}(\theta) \right) = \sqrt{1 - \theta^2}$.

regularity of the domain (which is a nonempty interval here) and the profile is known (and very smooth) within the interface and the interior of its complement²⁶. \diamond

Remark 47. Despite the fact that the analytical solution ϕ is C^1 and even C^∞ within the interface and its complements, the use of this monotonous transition profile requires “violating” Equation (6.26), and thus, by the complementary condition, re-introducing the multipliers μ^\pm . The only regions where ϕ need not satisfy this equation are the ones where $\phi = 0$ or $\phi = 1$, leading to $(\mu^-, \mu^+) = (\frac{16\gamma}{\pi^2\epsilon}, 0)$ (resp. the other way around). In addition, it is easy to see that it is not possible to match this inner solution in a C^2 -manner with the (constant) outer regions as - by the nature of the sin-function - the only points where $\frac{d^2\phi}{dx^2} = \frac{d^2\phi_h}{dx^2} = 0$ are the ones where the homogeneous solution itself vanishes and therefore $\phi(x) = \frac{1}{2} \notin \{0, 1\}$. Whereas the second derivative in the outside region is obviously zero, the limits from the interior are given by $\pm \frac{8}{\pi^2\epsilon^2}$, showing in particular that the phasefield equation with the obstacle potential above does not allow for a classical solution unless it is sinusoidal everywhere or if there is no interface at all (i.e. $\phi(x) \equiv 0$ or 1). \diamond

Remark 48. The jump in the second derivatives of course also has numerical consequences. Taking the analytical profile and applying e.g. a centered finite difference scheme for the second derivative at a point at the transition between the bulk and the interface, the discrete estimate for the (non-existent) second-order derivative at this point will converge to the average of the second derivatives on the left and right of the transition point²⁷ and thus, one of them being zero, to half the second derivative corresponding to the limit from the interior. This is of course not a problem per se as there not being a “correct” value to converge to at this point, the average is actually a quite pleasing limit. Nevertheless, a similar effect will happen for any point whose stencil crosses the purported transition point. Taking for definiteness the transition between a zero bulk-region and the interior, a point just outside the interface will overestimate the (vanishing) second derivative at this point, whereas one just on the inside will underestimate it, with the degree of both effects depending on the relative positions. Starting a discrete simulation from the analytical profile, one would therefore in particular expect the value of this

²⁶Assuming the actual domain Ω to be of the form $\Omega = (-R, R)$ with some sufficiently large R such that it fully contains the interface (the energetic contributions in the bulk vanishing), the first order necessary condition in its weak form is given by

$$\int_{-R}^R 2\gamma\epsilon \frac{d\phi}{dx} \frac{d\psi}{dx} + \frac{1}{\epsilon} w'(\phi)\psi \, dx = \int_{(-R, x_l) \cup (x_r, R)} \frac{1}{\epsilon} w'(\phi)\psi \, dx + \int_{x_l}^{x_r} 2\gamma\epsilon \frac{d\phi}{dx} \frac{d\psi}{dx} + \frac{1}{\epsilon} w'(\phi)\psi \, dx \geq 0$$

for all admissible directions ψ . Due to the smoothness of ϕ in the interface, the integration by part (corresponding to the use of Gauss’s theorem in the multidimensional case) leading to

$$\begin{aligned} & \int_{(-R, x_l) \cup (x_r, R)} \frac{1}{\epsilon} w'(\phi)\psi \, dx + 2\gamma\epsilon \left(\frac{d\phi}{dx}(x_r^-)\psi(x_r) - \frac{d\phi}{dx}(x_l^+)\psi(x_l) \right) - \int_{x_l}^{x_r} \left(-2\gamma\epsilon \frac{d^2\phi}{dx^2} + \frac{1}{\epsilon} w'(\phi) \right) \psi \, dx \\ &= \int_{(-R, x_l) \cup (x_r, R)} -\frac{16\gamma}{\pi^2\epsilon} \psi \, dx + 2\gamma\epsilon \left(\frac{d\phi}{dx}(x_r^-)\psi(x_r) - \frac{d\phi}{dx}(x_l^-)\psi(x_l) \right) \geq 0 \end{aligned}$$

in the interface is now easily justified provided ψ is moderately smooth and one takes the limits “from the interior” for the values of $\frac{d\phi}{dx}$. As there are no gradients involved in the expression for the bulk-region, one can then simply take e.g. a sufficiently narrow hat-function of the appropriate sign centered around each endpoint for ψ to conclude that $\frac{d\phi}{dx}(x_l)$ cannot be positive and $\frac{d\phi}{dx}(x_r)$ cannot be negative, contradicting the assumed shape of ϕ unless $\tilde{c}_1 = 1$.

²⁷This is easily seen by a two-sided Taylor-expansion, which, while not legitimate **through** the discontinuity of the second derivative at the transitions point, is valid in both the left and right region of the endpoints. Fixing the position x_i as one of two transition points x_l and x_r , a simple Taylor expansion using the left resp. right limit of the second derivative gives $\phi_{i-1} = \phi(x_i - \Delta x) = \phi(x_i) - \phi'(x_i)\Delta x + \frac{1}{2}\phi''(x_i^-)(\Delta x)^2 + \mathcal{O}((\Delta x)^3)$ and $\phi_{i+1} = \phi(x_i + \Delta x) = \phi(x_i) + \phi'(x_i)\Delta x + \frac{1}{2}\phi''(x_i^+)(\Delta x)^2 + \mathcal{O}((\Delta x)^3)$, from which the evaluation of the discrete (3-point stencil based) second derivative is obtained as $\frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} = \frac{1}{2}(\phi''(x_i^-) + \phi''(x_i^+)) + \mathcal{O}(\Delta x)$.

first inner point to start to drop²⁸. Similarly, the last inner point near $\phi = 1$ is expected to rise. This can serve to provide an intuitive explanation for the observation (see Figure 6.1 and the next section) that the resulting discrete phasefield profile is generally slightly more narrow and steeper than the continuous one if the last “inner” point is close to the transition to the bulk.

A further quite interesting observation (discussed in more detail in Appendix A) is that, at least without driving forces, there are generically (i.e. except for a discrete sequence of ratios of $\frac{\epsilon}{\Delta x}$) two possible solutions to the discrete first-order necessary conditions corresponding to Equation (6.22) with only one of them actually minimizing the discrete energy. For low resolutions, these differ quite notably as is also illustrated in Figure 6.1. There, N corresponds to the total distance in grid-spacings between the last discrete ϕ -value at 0 and the first one at 1, there thus being a total of $N - 1$ points within the interior of the discrete interface (i.e. with $0 < \phi_i < 1$) and a total of $N + 1$ when the two endpoints at 0 resp. 1 are included.

The discrete solution corresponding to $N = 4$ experiences precisely the effect just outlined above. As the difference stencil of the last discrete point at $\frac{x}{\epsilon} = 1$ (an integer here for the choice $\epsilon = 2$) is roughly at a distance of $\frac{\Delta x}{2}$ from the outer limit of the continuous profile, it will significantly underestimate the “true” curvature at this point and is thus forced upwards (into the upper bound at 1) as this weakens the contribution by the second derivative as compared to the bulk potential one.

In contrast, the solution corresponding to $N = 5$ is in a sense “numerically optimal” as the last discrete position is almost perfectly aligned with the right endpoint of the discrete profile (and only very slightly outside the actual interface indicated by the vertical dashed line). For this reason, the evaluation of the second derivative of the analytical solution at the last inner point is essentially unaffected by the jump discontinuity and thus quite precise (up to the standard $\mathcal{O}((\Delta x)^2)$ error, see Footnote 28). Even though the (negative) second derivative “just outside” the interface at the last point is highly overestimated, this only amounts (in line with the above) to roughly half of what is required to counterbalance the forcing by the bulk potential term. This leads to an underestimation of the corresponding multiplier μ^+ by again roughly $\frac{1}{2}$, but this is not sufficient to move the right-most point away from one. The jump of the second derivative is thus essentially “hidden” in the multipliers, explaining the very close numerical agreement for the profile itself. In this respect, it can also be noted that the error in the multiplier is much lower for the more narrow profile.

While this may therefore seem very favorable from an approximation point of view, the discussion in the next section and the Appendix A show that on energetic grounds the numerics should a priori favor the **other** profile with $N = 4$ as it has a (slightly) lower total “discrete” energy (of $\approx 0.981\gamma$ as compared to $\approx 0.984\gamma$ for $N = 5$). \diamond

²⁸The discrete second-order derivatives are accurate to second order at all points not in the neighborhood of the transition region and thus “almost” in equilibrium with the local $w'(\phi)$ -values. Even though a closer look reveals that the discrete three-point stencil applied to the analytical profile also underestimates the magnitude of the actual second derivative by a factor $2\frac{\cos(\Delta x)}{\Delta x} = 1 - \frac{2}{4!}(\Delta x)^2 + \mathcal{O}((\Delta x)^4)$, this is a priori nowhere near the reduction by a factor of up to 2 at the transition points. Nevertheless, this (perhaps) surprisingly turns out to have an effect of roughly the same order of magnitude (see Remark 66 for a more detailed discussion).

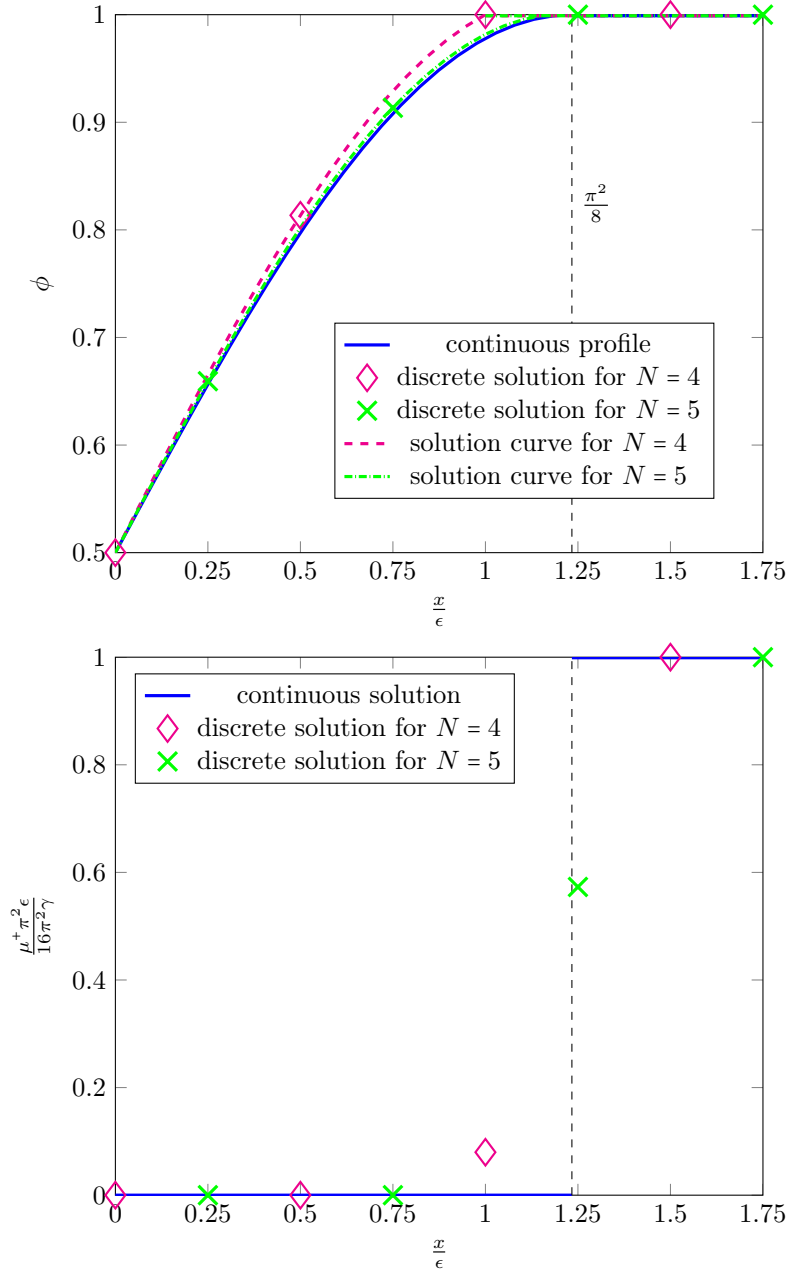


Figure 6.1: Comparison of the right half of the ϕ -profiles and their corresponding multiplier μ^+ for the continuous and the two (!) discrete solutions of the discrete analogue of the first-order necessary conditions in Equation (6.22) for $\epsilon = 2\Delta x$ (see Remark 48 and the next Subsection 6.2.3 for a precise definition of N).

For the discrete profiles, the dashed resp. dash-dotted lines in addition show the underlying sinusoidals on which the solution values lie (see Equation (6.48)). Note that while the curve for $N = 5$ is closer to the analytical profile, it is actually the energetically unfavorable one.

6.2.3 Some Analytics for the Discrete Case

The Discrete One-Dimensional Minimizers

A similar though somewhat heavier calculation for the obstacle potential can in fact, again relying on the linearity, be reproduced in the discrete setting and sheds some light on the differences between the continuous and discrete behaviour within the simple one-dimensional setting. Using a standard centered difference discretization for the Laplacian and discretizing w using the values at the cell centers, the discrete one-dimensional phasefield equation in the absence of driving forces is given by

$$-2\gamma\epsilon\frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon}\gamma(1 - 2\phi_i) = 0. \quad (6.31)$$

By linearity, one can again decompose the solution as the sum of the same particular solution $\phi_i^p = \frac{1}{2}$ as in the continuous case and the homogeneous solution ϕ_i^h satisfying

$$\phi_{i+1}^h - 2\left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)\phi_i^h + \phi_{i-1}^h = 0.$$

From the standard theory of difference equations (see e.g. [59]), the solution can be derived in terms of the solutions λ of the characteristic equation²⁹

$$\lambda^2 - 2\left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)\lambda + 1 = \left(\lambda - \left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)\right)^2 + \left(1 - \left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)\right)^2 = 0,$$

which, for $\left(\frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)^2 \leq 2$, are given by

$$\alpha + i\beta = \left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right) \pm i\sqrt{1 - \left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)^2}$$

where $i = \sqrt{-1}$ and $\alpha^2 + \beta^2 = 1$. The general homogeneous solution is therefore given by

$$\phi_i^h = c_1 \cos(\kappa i) + c_2 \sin(\kappa i)$$

with

$$\cos(\kappa) = 1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2} = 1 - \frac{1}{2}\left(\frac{4\Delta x}{\pi\epsilon}\right)^2 \quad (6.32)$$

$$\sin(\kappa) = \sqrt{1 - \left(1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)^2} = \sqrt{2\frac{8(\Delta x)^2}{\pi^2\epsilon^2} - \left(\frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)^2} = \frac{8(\Delta x)^2}{\pi^2\epsilon^2} \sqrt{\frac{\pi^2\epsilon^2}{4(\Delta x)^2} - 1}. \quad (6.33)$$

Due to the translation invariance of the problem, one may for simplicity assume that $\phi_0 = 0$, and that the difference equation holds from ϕ_1 up to ϕ_{N-1} . As this imposes $\phi_0^h = -\frac{1}{2}$, this fixes the first constant as $c_1 = -\frac{1}{2}$. Similarly, using the condition

$$\phi_N^h = -\frac{1}{2} \cos(\kappa N) + c_2 \sin(\kappa N) \stackrel{!}{=} \frac{1}{2} \quad (6.34)$$

²⁹The approach here is similar to the one in the continuous case. Instead of the ansatz $\phi_h(x) = e^{\lambda x}$ based upon which one can obtain the characteristic equation underlying the homogeneous solution in the continuous double-well case, one instead postulates $\phi_i = \lambda^i$ (with generally complex λ), from which one then obtains the characteristic equation by extracting the highest common power of λ . In the case with two complex-conjugate λ 's below, it is then convenient to rewrite λ^i as $e^{i \ln \lambda} = \cos(\Im(\ln \lambda)) + i \sin(\Im(\ln \lambda))$. Further rewriting λ in polar form as $\lambda = |\lambda|(\cos(\kappa) + i \sin(\kappa))$ and arguing precisely as in the continuous case then leads to the representation used here.

at the right endpoint leads to

$$c_2 = \frac{1 + \cos(\kappa N)}{2 \sin(\kappa N)} = \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \quad (6.35)$$

unless κN is a multiple of π and fixes the second constant in terms of the (as of yet unknown!) number of points N , leading to the discrete profile

$$\phi_i = \frac{1}{2} - \frac{1}{2} \cos(\kappa i) + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa i). \quad (6.36)$$

Remark 49. If κN is multiple of π , equation (6.34) does not impose any restriction on c_2 as $\cos(\kappa N) = 1$ and $\sin(\kappa N) = 0$, implying that this equation is satisfied for any c_2 . This particular case (for $\kappa N = \pi$) turns out to be a very interesting one as will be seen towards the end of this section. In order to avoid mixing two quite different lines of argument, the following paragraph will first consider the case where $\frac{\pi}{\kappa}$ is noninteger. \diamond

The case $\frac{\pi}{\kappa}$ noninteger All which remains to be done is therefore to determine the width of the numerical interface in terms of N . A first point to be noted is that the assumption of ϕ_1 being the first inner interface point already imposes an upper on N as this requires

$$\phi_1 = \frac{1}{2} - \frac{1}{2} \cos(\kappa) + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa) > 0.$$

Since $\sin(\kappa) > 0$, this translates to

$$\cot\left(\frac{\kappa N}{2}\right) > \frac{\cos(\kappa) - 1}{\sin(\kappa)} = -\tan\left(\frac{\kappa}{2}\right)$$

and thus, restricting the range of \cot^{-1} to the relevant one in $(0, \pi)$ to obtain a single-valued mapping and using $\cot^{-1}\left(\frac{1}{y}\right) = \frac{\pi}{2} - \cot^{-1}(y)$ for $y > 0$ and $\tan(y) = \left(\cot(y)\right)^{-1}$, to

$$N < \frac{2}{\kappa} \left(\pi - \cot^{-1} \left(\tan \left(\frac{\kappa}{2} \right) \right) \right) = \frac{2}{\kappa} \left(\pi - \left(\frac{\pi}{2} - \frac{\kappa}{2} \right) \right) = \frac{\pi}{\kappa} + 1.$$

Since N has to be integer, this together with the strict inequality leads to the maximal admissible number

$$N_{max} = \left\lceil \frac{\pi}{\kappa} \right\rceil \quad (6.37)$$

of points within the interface, where $\lceil \cdot \rceil$ denotes “round-up” operation (i.e. $\lceil n \rceil$ is the smallest integer number greater than or equal to n).

As all $N \leq N_{max}$ would in principle permit the construction of a “geometrically admissible” interface, an additional criterion is necessary for choosing the correct one. A first such criterion is easy to obtain by reading the multipliers μ^\pm corresponding to the constraint $0 \leq \phi_i$ and $\phi_i \leq 1$ ³⁰,

$$-2\gamma\epsilon \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon} \gamma (1 - 2\phi_i) = \mu_i^- - \mu_i^+, \quad (6.38)$$

with $\mu_i^\pm \geq 0$ and satisfying $\mu^- = 0$ if $\phi > 0$ and $\mu^+ = 0$ if $\phi < 1$, which extends the validity of the discrete difference Equation (6.31) to the entire domain. It is obvious that within the bulk-region one has either $\mu^- = \frac{16}{\pi^2\epsilon}$ or $\mu^+ = \frac{16}{\pi^2\epsilon}$, the respective other one being zero, and thus multipliers consistent with the restriction in terms of the complementarity condition (6.21). The only interesting points for judging the first-order admissibility for a given N are therefore the

³⁰Recall the reasoning for Equation (6.22), which automatically becomes rigorous in the discrete case.

transition points, where there is an additional contribution by the spatial second-order difference operator.

Inserting the profile in Equation (6.36) and looking for simplicity at the left end-point of the interface (the other endpoint being completely analogous) leads to the restriction

$$\mu_0^- = -2\gamma\epsilon \frac{\phi_1 - 2\phi_0 + \phi_{-1}}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon} \gamma(1 - 2\phi_0) = \gamma\epsilon \frac{-\cot\left(\frac{\kappa N}{2}\right) \sin(\kappa) + \cos(\kappa) - 1}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon} \gamma \stackrel{!}{\geq} 0 \quad (6.39)$$

since $\phi_{-1} = \phi_0 = 0$ and $\phi_1 = \frac{1}{2} - \frac{1}{2} \cos(\kappa) + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa)$. Canceling the common factor γ , it follows that N has to be such that

$$\cot\left(\frac{\kappa N}{2}\right) \leq \frac{1}{\sin(\kappa)} \left(\frac{16(\Delta x)^2}{\pi^2\epsilon^2} + \cos(\kappa) - 1 \right).$$

By the definition of κ in Equations (6.32) and (6.33), this can be simplified to

$$\cot\left(\frac{\kappa N}{2}\right) \leq \frac{1}{\sin(\kappa)} \left(\frac{16(\Delta \bar{x})^2}{\pi^2} - \frac{8(\Delta \bar{x})^2}{\pi^2} \right) = \frac{8(\Delta \bar{x})^2}{\pi^2 \sin(\kappa)} = \frac{1}{\sqrt{\frac{\pi^2}{4(\Delta \bar{x})^2} - 1}}$$

which represents a lower bound on N as the cot-function is decreasing from $+\infty$ at 0 to $-\infty$ at π (this being the interesting range here). More precisely, N has to satisfy $\frac{\kappa N}{2} \geq \cot^{-1}\left(\frac{1}{\sqrt{\frac{\pi^2}{4(\Delta \bar{x})^2} - 1}}\right)$, or, using the inverse formula $\cot^{-1}\left(\frac{1}{y}\right) = \tan^{-1}(y)$ and as N has to be integer,

$$N \geq N_{min} := \left\lceil \frac{2}{\kappa} \tan^{-1}\left(\sqrt{\frac{\pi^2\epsilon^2}{4(\Delta x)^2} - 1}\right) \right\rceil = \left\lceil \frac{\pi}{\kappa} - 1 \right\rceil \quad (6.40)$$

Any (integer) N satisfying both equations (6.37) and (6.40) is therefore compatible with the first-order optimality conditions in Equation (6.38). As it is obvious from the expressions for N_{min} and N_{max} that there are always two such integers, this shows that the existence of two potentially relevant profiles in the example in Figure 6.1 is not an exception but actually the rule.

Since the discrete energy is not convex due to the concave contribution by w and this condition is in the non-convex case only necessary but **not** sufficient for characterizing a local minimizer, this naturally raises the question as to the nature of the critical points attained by both profiles. A first insight into this matter can be obtained by comparing the discrete energies

$$\mathcal{E}(N) = \epsilon \sum_{i=0}^N \frac{1}{2} \gamma \left(\left(\frac{\phi_i - \phi_{i-1}}{\Delta x} \right)^2 + \left(\frac{\phi_{i+1} - \phi_i}{\Delta x} \right)^2 \right) + \frac{16\gamma}{\pi^2\epsilon} \sum_{i=0}^N \phi_i (1 - \phi_i) \quad (6.41)$$

for both potential profiles (this being the one for which Equation (6.38) is the FONC).

Remark 50. This is similar to the discussion in Remark 46, except that guessing the correct choice is more difficult in the discrete case as enforcing a continuous derivative is meaningless in this setting. \diamond

Remark 51. Note that in order to obtain a discrete approximation of the continuous energy, this expression still needs to be multiplied by Δx to take the physical size of each cell into account. The first contribution to this approximation is then given by $\frac{4\gamma}{\pi^2\epsilon} N \Delta x$. As Δx converges to zero, $N \Delta x$ (as the width of the discrete interface) will converge to the width of the analytical one in Equation (6.30) (i.e. to $\frac{\pi^2\epsilon}{4}$) and the first term therefore to the desired value of γ . It further follows that the second has to (and actually does) converge to zero.

A more quantitative evaluation of the energetic deviations requires considering both terms as $N \Delta x$ will only reduce to the correct value in the limit. \diamond

Even though this is a rather lengthy calculation (see the Appendix A), the discrete energy contributions in Equation (6.41) can be evaluated analytically, leading to

$$\sum_{i=1}^{N-1} \frac{1}{\epsilon} w(\phi_i) = \frac{8}{\pi^2 \epsilon} \gamma \left(\frac{1}{4} \left(1 - \cot^2 \left(\frac{\kappa N}{2} \right) \right) N + \frac{1}{2} \cot \left(\frac{\kappa N}{2} \right) \cot(\kappa) \right) \quad (6.42)$$

and

$$\sum_{i=1}^{N-1} \epsilon a_i = \frac{8}{\pi^2 \epsilon} \gamma \left(\frac{1}{4} \left(1 + \cot^2 \left(\frac{\kappa N}{2} \right) \right) N + \frac{1}{2} \frac{\cot \left(\frac{\kappa N}{2} \right)}{\sin(\kappa)} \right) = \frac{8}{\pi^2 \epsilon} \gamma \left(\frac{N}{4 \sin^2 \left(\frac{\kappa N}{2} \right)} + \frac{1}{2} \frac{\cot \left(\frac{\kappa N}{2} \right)}{\sin(\kappa)} \right), \quad (6.43)$$

or, using $1 - \cot^2(\theta) = 2 - \frac{1}{\sin^2(\theta)}$ and $1 + \cot^2(\theta) = \frac{1}{\sin^2(\theta)}$, alternatively to

$$\sum_{i=1}^{N-1} \frac{1}{\epsilon} w(\phi_i) = \frac{8}{\pi^2 \epsilon} \gamma \left(\left(\frac{1}{2} - \frac{1}{4 \sin^2 \left(\frac{\kappa N}{2} \right)} \right) N + \frac{1}{2} \cot \left(\frac{\kappa N}{2} \right) \cot(\kappa) \right) \quad (6.44)$$

and

$$\sum_{i=1}^{N-1} \epsilon a_i = \frac{8}{\pi^2 \epsilon} \gamma \left(\frac{1}{4 \sin^2 \left(\frac{\kappa N}{2} \right)} N + \frac{1}{2} \frac{\cot \left(\frac{\kappa N}{2} \right)}{\sin(\kappa)} \right). \quad (6.45)$$

Combining both expressions and making use of the defining properties of κ (see Appendix A for details), this can in addition be shown to lead to a total energy given by

$$\mathcal{E}(N) = \frac{4\gamma}{\pi^2 \epsilon} N + \frac{\epsilon \gamma}{(\Delta x)^2} \frac{1}{2} \cot \left(\frac{\kappa N}{2} \right) \sin(\kappa) = \frac{4\gamma}{\pi^2 \epsilon} \left(N + \cot \left(\frac{\kappa N}{2} \right) \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1} \right). \quad (6.46)$$

A further analysis (see again Appendix A) then shows that, from an energetic point of view, unless $\frac{\pi}{\kappa}$ happens to be an integer, the energetically optimal interface width is given by the lower bound on N in Equation (6.40), i.e. whenever there is more than one potentially admissible discrete interface width, one has $N_{opt} = N_{min}$, leading to the optimal profile being determined by

$$N = N_{min} = \left\lceil \frac{\pi}{\kappa} - 1 \right\rceil. \quad (6.47)$$

From this, the profile itself and the associated discrete energies follow immediately from equations (6.36), (6.42) and (6.43) (resp. (6.44) and (6.45)) as well as (6.46) using the actual number $N = N_{min}$ of points.

Remark 52. The discrete profile is, as in the continuous case, only fixed up to a reflection and an arbitrary translation as long as this does not interfere with the boundaries of the domain. Even though the choice of fixing the profile as an increasing one starting at $i = 0$ is convenient for some of the calculations above and in Appendix A, one can of course also recenter the profile at the 0.5-isoline, i.e. the point $x_m = \frac{N}{2} \Delta x$. Using

$$\cos(\kappa i) = \cos(\kappa(x_i - x_m) + \kappa x_m) = \cos(\kappa(x_i - x_m)) \cos(\kappa x_m) - \sin(\kappa(x_i - x_m)) \sin(\kappa x_m)$$

and

$$\sin(\kappa i) = \sin(\kappa(x_i - x_m) + \kappa x_m) = \sin(\kappa(x_i - x_m)) \cos(\kappa x_m) + \cos(\kappa(x_i - x_m)) \sin(\kappa x_m),$$

leads to

$$\begin{aligned} \phi_i &= \frac{1}{2} - \frac{1}{2} \cos(\kappa i) + \frac{1}{2} \cot \left(\frac{\kappa}{2} N \right) \sin(\kappa i) \\ &= \frac{1}{2} - \frac{1}{2} \left(\cos(\kappa x_m) - \cot \left(\frac{\kappa}{2} N \right) \sin(\kappa x_m) \right) \cos(\kappa(x_i - x_m)) \\ &\quad + \frac{1}{2} \left(\sin(\kappa x_m) + \cot \left(\frac{\kappa}{2} N \right) \cos(\kappa x_m) \right) \sin(\kappa(x_i - x_m)). \end{aligned}$$

As, by the definition of x_m , $\kappa x_m = \frac{\kappa N}{2}$ and $\cos\left(\frac{\kappa}{2}N\right) - \cot\left(\frac{\kappa}{2}N\right)\sin\left(\frac{\kappa}{2}N\right) = 0$, the prefactor for the cosine cancels³¹, whereas the one for the sine-term simplifies to

$$\sin\left(\frac{\kappa}{2}N\right) + \cot\left(\frac{\kappa}{2}N\right)\cos\left(\frac{\kappa}{2}N\right) = \sin\left(\frac{\kappa}{2}N\right) + \frac{\cos^2\left(\frac{\kappa}{2}N\right)}{\sin\left(\frac{\kappa}{2}N\right)} = \frac{1}{\sin\left(\frac{\kappa}{2}N\right)},$$

leading to the more pleasant representation

$$\phi_i = \frac{1}{2} + \frac{1}{2} \frac{1}{\sin\left(\frac{\kappa}{2}N\right)} \sin\left(\frac{\kappa}{\Delta x}(x_i - x_m)\right). \quad (6.48)$$

The other choice of orientation (i.e. a decreasing profile) is then obviously again obtained by simply changing the sign of the second term. \diamond

In contrast to the choice for $N = N_{min}$, it turns out that the second possible choice of $N = N_{max} = N_{min} + 1$ does **not** satisfy the second-order necessary conditions for the characterization of a local minimizer in terms of the phasefield, i.e. this second choice is not only not optimal in terms of N , but does also not correspond to a local minimizer (of which there might a priori be more than one) but simply to a critical point of the functional \mathcal{E} .

The key steps for obtaining this result are (see Appendix A for additional details) that, since the multipliers μ^\pm are, for this larger choice of N , strictly positive, the second-order necessary condition in fact reduces to the study of the local stability for a **given** N , and thus, by the linearity of the equations **within** the interface, to showing the positive definiteness of the system matrix corresponding to the difference stencil in Equation (6.31). As this is simply a difference operator with constant coefficients (meaning also that the stability in particular has nothing to do with the ϕ -values themselves, but only with the choice of N), this reduces to an examination of the lowest eigenvalue of the discrete second-order difference operator under vanishing Dirichlet boundary data.

Remark 53. As the analysis below will show, κ is, up to a third order error in $\frac{\Delta x}{\epsilon}$, given by $\frac{4\Delta x}{\pi\epsilon}$, whereas N is, up to a term of order $\mathcal{O}(1)$ (linked to the discrepancy between N_{min} and N_{max} in equations (6.40) and (6.37)) given by $\frac{\pi^2\epsilon}{4\Delta x}$, i.e. the “continuous” number of points necessary to subdivide the analytical interface width into cells of width Δx . The critical dependence of the stability on the question of whether κN is less than or greater than π is therefore quite intuitive as this essentially measures whether the discrete interface is broader or more narrow than the analytical one. Roughly speaking, the choices of N for which the last point is slightly within the the continuous interface (i.e. the ones with the more narrow profiles) are stable, whereas those for which the last point is at the end or outside the continuous interface (i.e. the slightly broader profiles) are unstable (see the previous Figure 6.1 and Figure 6.4, which shows one of the “critical” points). \diamond

Remark 54. Even though the choice with $N = N_{min} + 1$ is not a stable solution for a descent-based algorithm, it can sometimes nevertheless be obtained (and in some cases actually easier so than the actual solution) in numerical simulations. This somewhat surprising observation is likely due to the fact that all eigenvectors (for this fixed N) which are **antisymmetric** with respect to the center of the interface (i.e. maintain the basic symmetry of the profile) are indeed stable ones. Even without actual stability, this “stability on symmetric profiles” can persist for quite some time when evaluated based on a numerically symmetric calculation³². \diamond

³¹This also being clear by construction as the value at $x_i = x_m$ has to be 0.5, which can only happen if the cosine is eliminated.

³²One example for this is in fact the example in Figure 6.1 for $\epsilon = 2\Delta x$. Starting from a sharp transition, a basic gradient-descent scheme tends to develop into the profile with $N = 5$ instead of the predicted value $N = 4$ and remains stable there (at least for a long time). This profile can be “broken down” though e.g. by perturbing the ϕ -values at the first inner point by a term of order 10^{-11} , which, with some delay, then leads to the predicted profile.

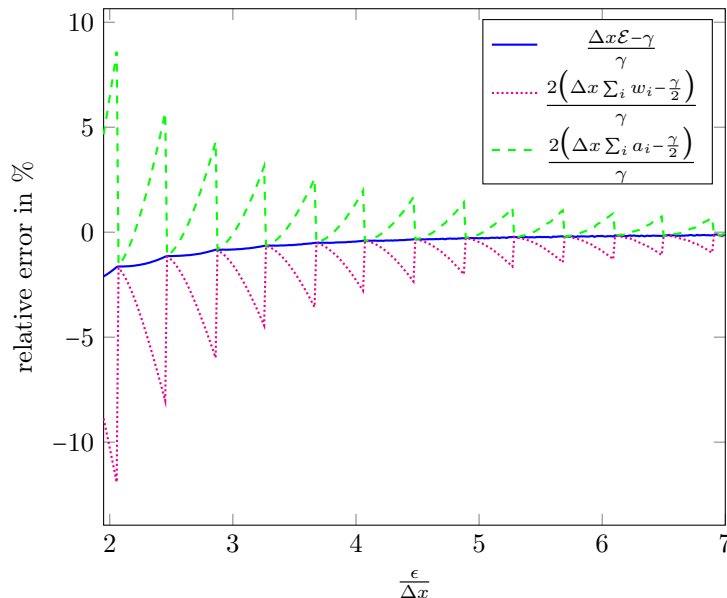


Figure 6.2: Development of the relative errors in the total energy and the bulk potential and gradient energy as a function of the ratio $\frac{\epsilon}{\Delta x}$ (normalized w.r.t the exact values γ for the total energy resp. $\frac{\gamma}{2}$ for the individual contributions).

One interesting consequence which can be derived from the explicit expressions for the total energy and its two contributions is that it allows for a straightforward analysis of the energy landscape and how it affects the discrete interface width³³. For a more intuitive approach, this energy landscape (in terms of the relative developments of the total energy $\Delta x \mathcal{E}$ and that of the contributions by the a - and w -term) is plotted in Figure 6.2.

A first observation to be made is that the total energy is in fact quite accurate and develops relatively smoothly, even for very low resolutions³⁴. The contributions by the the bulk potential and gradient energy in contrast exhibit a strongly oscillating and discontinuous behavior, with significantly larger deviations than the total energy. The accuracy of the latter is therefore primarily the result of cancellations, as is also obvious from Figure 6.2.

This - as well as the tendency of the signs of the deviations - is intuitively relatively easy to understand based on the competing nature of a and w combined with Figure 6.3, which shows the relative errors in the interface width (counted as the distance between the last value at 0 and the first one at 1). There, it is obvious that the discrete interface is always - and for low resolutions quite significantly so - more narrow than the continuous one. As the role of the gradient energy term is to try to enlarge the interface width whereas that of the bulk energy potential is to make it thinner, it is clear by the decreased width that the discrete case consistently favors the bulk potential term. This leads to a deviation from the equipartition of energy in the continuous case, with the bulk potential energy being below its continuous value $\frac{\gamma}{2}$, but partially balanced by an increase in the gradient contribution due to the necessity of a steeper interface (also recall Figure 6.1)³⁵.

³³This could of course alternatively be evaluated based on a significantly more tedious simulation study.

³⁴Note that the left-most point has a relative deviation of less than two percent, even though the number of points (including the respectively first “bulk-point” on each side) is only 4 as already seen in Figure 6.1 and also shown in Figure 6.3.

³⁵The points where the gradient energy contribution drops below its reference value of $\frac{\gamma}{2}$ are those where the

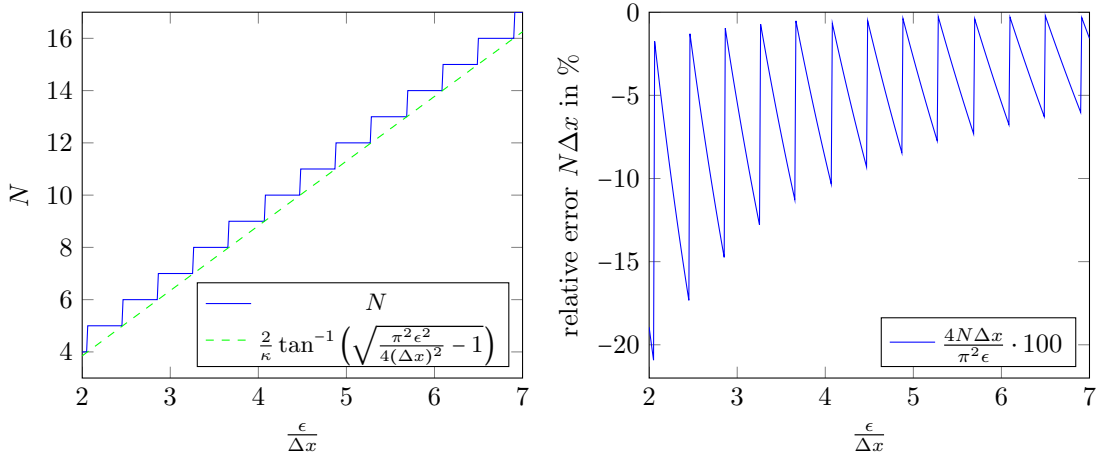


Figure 6.3: Development of the number of interface points and the relative error in the “discrete interface width” with respect to $\frac{\pi^2}{4}\epsilon$ (also see Remark 59).

Remark 55. The discontinuous nature of the problem is also evident from both figures. In particular, as the development of N in Figure 6.3 shows, the evolution of the number of interface points would almost be linear (as one would expect by simply subdividing the analytical interface width based on a spacing of Δx), but, by the necessity of the $[\cdot]$ -operation in Equation (6.47), experiences jump discontinuities.

Each of these jumps is accompanied by a sudden drop in the gradient energy contribution and a sudden increase in the bulk potential energy as seen in Figure 6.2. In contrast, the total energy remains relatively smooth as the jumps in both contributions seem to cancel each other. \diamond

Remark 56. It can also be noted that the profile shown in Figure 6.1 at the end of the previous Section 6.2.2 is, with $\frac{\epsilon}{\Delta x} = 2$, relatively close to one of the “critical” ratios associated with adding an additional discrete point in the interface. It is interesting to take a closer look at what happens just to the right of one of the discontinuity points visible in Figures 6.2 and 6.3.

The effect on the discrete ϕ -profile is exemplarily shown in Figure 6.4 for a ratio of $\frac{\epsilon}{\Delta x}$ just to the right of the discontinuity point near 2, with, for comparison, the profile for a nearly double ration of $\frac{\epsilon}{\Delta x} = 4$. As can be seen there, even though the profile is quite accurate for the higher number of interface points (with $N = 9$ i.e. a total number of 8 points strictly within the interface), the match for the significantly lower ratio is also extremely good (and in fact seemingly even slightly better) despite only having $N = 5$ and thus 4 “interior” interface points. While the much better match as compared to the ratio of 2 in Figure 6.1 is not surprising given the discussion above (and the fact that there is now one additional point), it should be kept in mind that there are two factors influencing the profile. The first one, given by the discrete angular frequency κ , is completely independent of the number of actual points within the interface and is (with roughly 2% compared to 0.5% deviation from the “perfect” value $\frac{\pi\Delta x}{4\epsilon}$, see the next Subsection 6.2.3) notably more accurate for the higher ratio $\frac{\epsilon}{\Delta x} = 4$. In contrast, for $\epsilon \approx 2.061\Delta x$, the value of κ is very close to the one of the “bad” profile for $\epsilon = 2\Delta x$. The

interface width is quite accurate, i.e. only slightly more narrow than the continuous one. The solution profile at these points turns out to be quite accurate (see Remark 56). As the evaluation of the discrete gradient corresponds to evaluating an “average” gradient over the width Δx , it is to be expected that the discrete gradient energy will be slightly lower than the continuous one since this corresponds to squaring the average gradient instead of “averaging” (through the continuous integration) the square, and one always has $(\frac{a+b}{2})^2 \leq \frac{1}{2}(a^2 + b^2)$ with strict inequality unless $a = b$.

second one, depending both on κ and N , is given by the amplitude $\sin\left(\frac{\kappa N}{2}\right)^{-1}$ of the sinusoidal in Equation (6.48). It is only through a partial cancellation of both influences due to the need of fitting a slightly too “fast” sinusoidal into an almost correct interface width that one obtains the observed excellent match.

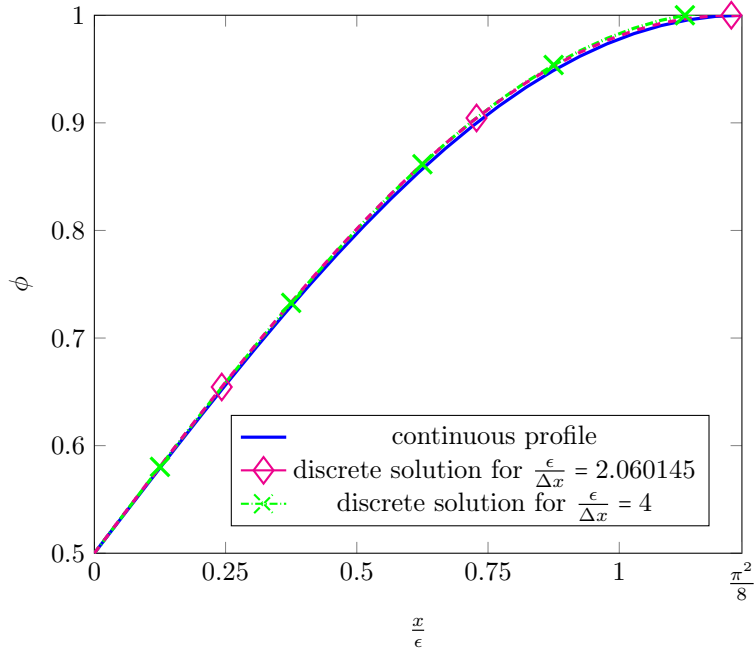


Figure 6.4: Discrete solutions together with the underlying sinusoids for $\frac{\epsilon}{\Delta x}$ just to the right of the critical ratio 2.06014486 (with $N = 5$) and for $\frac{\epsilon}{\Delta x} = 4$ (with $N = 9$) in comparison with the continuous profile. Note that even though the resolution in the second case is almost twice as high, the deviations in the profile itself are actually slightly larger than for the one just to the right of the critical point near $\frac{\epsilon}{\Delta x} = 2$.

A further interesting observation is that, as already discussed in Remark 48, the error in the multipliers μ^\pm is particularly high just to the right of the critical points as this will lead to an underestimation of the second derivative at the last inner point by a factor of almost 2. Even though one might - wrongly as will be explained below - conjecture that this could lead to problems related with the undesirable pinning of the discrete interface (i.e. its inability to move despite the action of a driving force) due to the large forcing required to move this point out of the 0 or 1 bound, this seems not to be the case. Instead, as the simple numerical test-case in Figure 6.5 illustrates, these regions at the transition between two different values of N in fact also seem to be “sweet-spots” for the mobility of the interface. The figure shows the final position of an interface initially placed near the left end of the domain under the influence of a (relatively weak) driving force pulling it towards the center of the domain. As is easily visible, the “arbitrary” low-resolution profile with $\epsilon = 2\Delta x$ leads to a relatively high deviation, whereas the one with $\epsilon = 4\Delta x$ (and thus the higher numerical resolution) as expected performs significantly better.

What is a priori somewhat unexpected is that, despite its markedly lower resolution of the interface, the profile with $\Delta x = 2.061\Delta x$ (a point just slightly to the right of the jump in Figure 6.2) not only results in a very accurate profile in the absence of driving forces as seen in Figure 6.4, but also leads to a final interface position which is much more accurate than the higher-resolution

one for $\epsilon = 4\Delta x$.

It can also be observed that a similar improvement in accuracy occurs for $\epsilon = 3.667\Delta x$, corresponding to a ratio just to the right of the jump to the left of $\frac{\epsilon}{\Delta x} = 4$ in Figure 6.2, which has the same number of “interface points” as the simulation for a ratio of four. Finally, the results for $\epsilon = 3.57\Delta x$ and $\epsilon = 3.77\Delta x$ - very close to being at the same distance below and above the critical value at $\frac{\epsilon}{\Delta x} \approx 3.66615$ - illustrate that this effect is quite sensitive to the choice of the ratio, but seems to favor slightly larger over slightly lower ratios.

This example illustrates that, even though the analysis above itself only covers the pure unperturbed one-dimensional case, the critical ratios at the transition points seem to maintain some interesting properties even in the presence of slight perturbations. As these are the values where $\kappa = \frac{\pi}{n}$, Equation (6.32) shows that these ratios are given by $\cos(\kappa) = \cos\left(\frac{\pi}{n}\right) = 1 - \frac{8(\Delta x)^2}{\pi^2 \epsilon^2}$ resp. after solving for $\frac{\epsilon}{\Delta x}$ using $1 - \cos(\theta) = 2 \sin\left(\frac{\theta}{2}\right)$ by

$$\frac{\epsilon}{\Delta x} = \frac{2}{\pi \sin\left(\frac{\pi}{2n}\right)}. \quad (6.49)$$

Some of these values are listed in Table 6.1.

As a moderate adjustment of a given ratio to the critical one just below it is something which can essentially be done for “free” (i.e. one expects to maintain the same number of “active” points for which actual calculations have to be performed), the vicinity of these ratios may therefore also have some practical interest for finite-difference based simulations in more complex settings, provided the interfaces are primarily aligned with the axis of the computational domain such as in e.g. some directional solidification or capillary rise problems.

In contrast, it is unlikely that they have any particular effect for settings with roughly spherically-shaped bulk regions, since the effective “grid-spacing” for the profiles along any oblique line can differ very significantly from the one along the axis (i.e. for example a factor of $\sqrt{2}$ in the two-dimensional case when an interface is oriented at an angle of 45°).

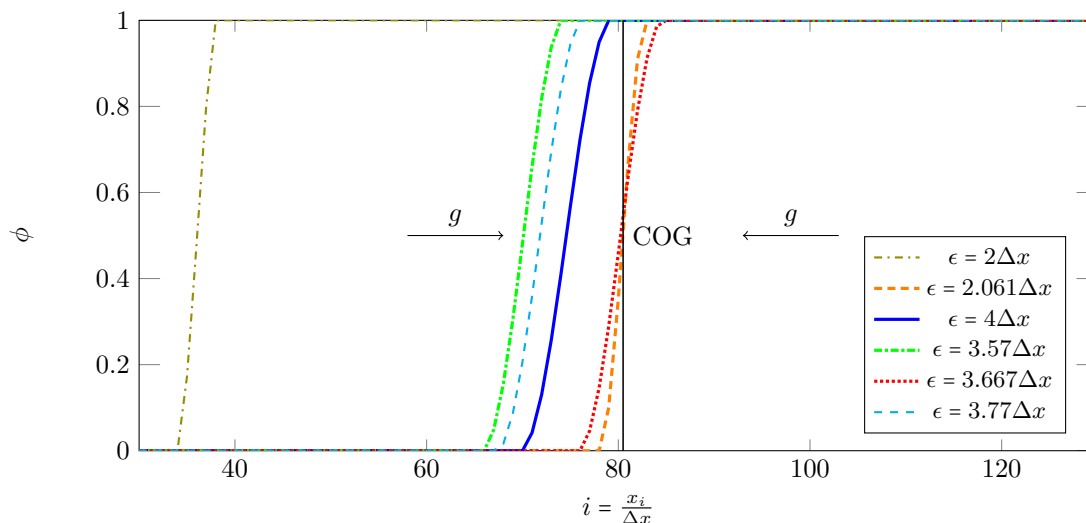


Figure 6.5: Final position of an interface initially placed on the left side under the influence of a “gravity” term pulling it towards the center of the gravity (COG) when varying ϵ for a fixed grid spacing Δx .

◇

N	5	6	7	8	9	10	11	12	13
$\left(\frac{\epsilon}{\Delta x}\right)_{min}$	2.460	2.861	3.264	3.667	4.070	4.474	4.878	5.282	5.686

Table 6.1: Some examples of the minimal values for the ratio of ϵ and Δx given by Equation (6.49) for a given number of points (rounded up at the fourth digit).

The case $\frac{\pi}{\kappa}$ integer As already alluded to in Remark 49 and seen on the example in Figure 6.5, the critical ratios for which $\frac{\pi}{\kappa}$ is integer are very particular and quite interesting. Unfortunately, much of the analysis from the previous paragraph does not carry over to the case where $\frac{\pi}{\kappa}$ is integer and requires different and slightly more complex arguments. For this reason, the discussion of this case is primarily performed in Appendix A, but the results will be summarized here.

Firstly, Equation (6.34) leaves c_2 open, which can be a priori be chosen arbitrarily as long as this is compatible with the bound constraint on the ϕ_i , $1 \leq i \leq N$. As the definition of c_2 in Equation (6.35) does not apply then, one cannot derive the upper bound $N_{max} = \lceil \frac{\pi}{\kappa} \rceil = \frac{\pi}{\kappa}$ based on the value of ϕ_1 (or ϕ_{N-1}) as above. Given that N does not explicitly enter the solution through c_2 , it is neither possible to derive the lower bound for N_{min} in Equation (6.40) based on the positivity of the multiplier μ_0^- (resp. the analogous estimate for μ_N^+ .)

Nevertheless, it turns out that both the bounds in terms of N_{min} and N_{max} continue to hold due to energetic considerations. Furthermore, choosing c_2 as in Equation (6.35), not only do the discrete energies for $N = N_{min}$ and $N = N_{max}$ in Equation (6.36) have the same energy, but there is actually a whole family of local minimizers with the same energy connecting these two “extreme” profiles. In fact, for $N = N_{max} = \frac{\pi}{\kappa}$, the sinusoidals with angular frequency κ are eigenfunctions of the homogeneous difference operator (i.e. dropping the constant $\frac{16}{\pi^2 \epsilon}$ arising from the w -term) within the interior of the interface in Equation (6.31) with an eigenvalue of 0 subject to vanishing boundary conditions at $i = 0$ and $i = N_{max}$. As a consequence, any multiple of this function can be added to the solution corresponding to N_{max} for $c_2 = \cot\left(\frac{\kappa N_{max}}{2}\right)$, which corresponds to the symmetric solution for this number of points, without affecting the validity of Equation (6.31). In addition, while a change of c_2 to $c_2 + \delta c_2$ will affect the values of μ_0^- and by symmetry μ_N^- in Equation (6.39) by $2\gamma\epsilon\delta c_2 \frac{\sin(\kappa)}{(\Delta x)^2}$ for μ_0^- resp. with the opposite sign for μ_N^+ , the common of the initial multiplier is given by $\frac{8}{\pi^2 \epsilon} \gamma$ (i.e. precisely half the bulk-value in the bulk) and is therefore strictly positive. As long as $|\delta c_2| \leq \frac{4(\Delta x)^2}{\pi^2 \epsilon \sin(\kappa)} = \frac{1}{2} \tan\left(\frac{\kappa}{2}\right)$, both multipliers at $i = 0$ and $i = N$ maintain the correct sign. Finally, the values $\delta c_2 = \pm \frac{1}{2} \tan\left(\frac{\kappa}{2}\right)$ are also the ones where either ϕ_{N-1} just touches 1 resp. where ϕ_1 just touches 0, i.e. the values where a new constraint becomes active corresponding to the solution (or a translate by one cell) for the value $N = N_{max} - 1$. More precisely, as

$$\frac{1}{2} \tan\left(\frac{\kappa}{2}\right) = \frac{1}{2} \cot\left(\frac{\pi}{2} - \frac{\kappa}{2}\right) = \frac{1}{2} \cot\left(\frac{\kappa\left(\frac{\pi}{\kappa} - 1\right)}{2}\right) = \frac{1}{2} \cot\left(\frac{\kappa N_{min}}{2}\right),$$

and $\cot\left(\frac{\kappa N_{max}}{2}\right) = 0$, the family of functions

$$\phi_i = \frac{1}{2} - \frac{1}{2} \cos(\kappa i) + c_2 \sin(\kappa i) \quad , 0 \leq c_2 \leq \frac{1}{2} \cot\left(\frac{\kappa N_{min}}{2}\right)$$

satisfies the first-order necessary conditions for a local minimizer, has a constant total energy and continuously transforms the “standard” profile with $c_2(N_{max})$ to the one for $c_2(N_{min})$.

This shows that there is an energetically neutral way of moving the right end of the broader interface one point to the left for obtaining the thinner interface and similarly, reverting the order of the argument for moving the left endpoint for the thinner interface one point to the left for again obtaining the broader interface, but now translated by one cell. A symmetric

discussion clearly also applies for moving the interface to the right, and then, by repetition, shows that there is an entire family of interfaces through which can move the interface, albeit with an “oscillating” deformation along the x -axis.

Remark 57. While this argument only strictly holds for the unperturbed equation at the precise critical ratios, it provides an intuitive explanation for the situation observed in Figure 6.5 since the “good” ratios are very close to a value where the interface is actually free to move without having to pass through any energetic barriers. \diamond

Remark 58. Even though the total energy is constant during the transformation above, this is not true for the gradient- and bulk-contributions separately. \diamond

An Asymptotic Analysis of the Profiles

Even though the expressions for the phasefield profile and discrete energy in equations (6.36) and (6.46) are exact and can easily be evaluated numerically, they are not very “readable” in this form. It is relatively straightforward to derive more tractable approximations though as Δx tends to zero.

One a priori problematic point here is that the definition $N = N_{true}$ in Equation (6.47) relies on a “round-up” operation. The number of points can therefore only be estimated up to a error of order 1, i.e.

$$N = \frac{\pi}{\kappa} - 1 + \delta, \quad \delta \in [0, 1] \quad (6.50)$$

where $\Delta \bar{x} := \frac{\Delta x}{\epsilon}$ abbreviates the ratio of the grid-spacing to the length-scale ϵ of the interface. It turns out that, disregarding this round-up, the width $N\Delta x$ of the discrete interface is, up to a second order term, smaller than the continuous value $\frac{\pi^2}{4}\epsilon$ by a single grid-spacing Δx . As this is a term of the same magnitude, the round-up, even though hard to estimate, is therefore not really limiting the precision of the general first-order convergence estimate for the interface width.

Considering first the value of κ as defined in Equation (6.32), an estimate for κ can be derived based upon the expansion

$$\cos(\kappa) = 1 - \frac{1}{2}\kappa^2 + \frac{1}{4!}\kappa^4 + \mathcal{O}(\kappa^6) \stackrel{!}{=} 1 - \frac{1}{2}\left(\frac{4}{\pi}\Delta \bar{x}\right)^2$$

of the cosine around zero (this being the relevant range for κ). An estimate based upon comparing the quadratic term only then shows that $\kappa = \mathcal{O}(\Delta \bar{x})$ and then more precisely with $\frac{\kappa}{\Delta \bar{x}} = \mathcal{O}(1)$ and $\sqrt{1-y^2} = 1 + \frac{1}{2}y^2 - \mathcal{O}(y^4)$ that³⁶

$$\kappa = \sqrt{\left(\frac{4}{\pi}\Delta \bar{x}\right)^2 + \mathcal{O}(\kappa^4)} = \frac{4}{\pi}\Delta \bar{x}\sqrt{1 + \mathcal{O}\left(\frac{\kappa^4}{(\Delta \bar{x})^2}\right)} = \frac{4}{\pi}\Delta \bar{x}\left(1 + \mathcal{O}\left((\Delta \bar{x})^2\right)\right). \quad (6.51)$$

Combined with $(1+y)^{-1} = 1 - y + \mathcal{O}(y^2)$, it follows that

$$\frac{1}{\kappa} = \frac{1}{\frac{4}{\pi}\Delta \bar{x}\left(1 + \mathcal{O}\left((\Delta \bar{x})^2\right)\right)} = \frac{\pi}{4\Delta \bar{x}} \frac{1}{1 + \mathcal{O}\left((\Delta \bar{x})^2\right)} = \frac{\pi}{4\Delta \bar{x}} \left(1 - \mathcal{O}\left((\Delta \bar{x})^2\right)\right),$$

from which the formula (6.50) for N can be rewritten as

$$N = \frac{\pi^2}{4\Delta \bar{x}} \left(1 - \mathcal{O}\left((\Delta \bar{x})^2\right)\right) - 1 + \delta = \frac{\pi^2}{4\Delta \bar{x}} - 1 + \delta - \mathcal{O}(\Delta \bar{x}), \quad \delta \in [0, 1]. \quad (6.52)$$

³⁶A more precise estimate will be derived in Equation (6.57).

Remark 59. Note that by the construction above, $\phi_0 \sim \phi(x_0) = \phi(0)$ is the last point where the ϕ -value is still 0, and $\phi_N \sim \phi(x_N) = \phi(N\Delta x)$ is the first point where the ϕ -value is at 1. Even though there is not really a precise definition of the interface width in the discrete case³⁷, the value of $N\Delta x = \frac{\pi^2}{4}\epsilon - (1 - \delta)\Delta x$ can be considered to be an upper bound for any reasonable interpretation of the discrete interface width, and the discrete interface is therefore slightly more narrow (with $-1 + \delta < 0$) than the continuous one. \diamond

Remark 60. Note that the approximation for κ shows that³⁸ the angular frequency $\frac{4}{\pi\epsilon}$ of the continuous solution in Equation (6.30) is achieved up to a second-order error.

Even though the estimate for $\frac{1}{\kappa}$ is only an estimate up to the order $\mathcal{O}(\Delta\bar{x})$, it is so around a value tending to infinity. Multiplying N by Δx to obtain the physical length of the interface, this error (with the bounded \tan^{-1} -term) only leads to a second-order deviation, whereas the value $N\Delta x$ is generally only first-order accurate.

What is potentially surprising is that, based on Equation (6.52), the a priori fully numerically induced necessity of the round-up operation turns out to actually cancel part of this first-order error and that, if δ is very close to 1, the interface will match, then essentially up to second order, the continuous one. This is also consistent with Figure 6.4, where the remarkable agreement of the discrete profile with the continuous one, despite the low resolution, is primarily due to the fact that the “continuous” approximation for N is just slightly above 4, and therefore, together with the $\lceil \cdot \rceil$ -operation leads to $\delta \approx 1$.

That the first-order error in the approximate width of the domain does not entail an error of equal order in the solution can intuitively be explained by the fact that the continuous solution enters the bulk with a slope of 0. While this is an unpleasant property in the “forward” direction (i.e. the ratio of the change of the interface width as compared to a change in the solution can be made arbitrarily high), it equally well applies in the “backward” direction, i.e. one can make the change in the profile despite a perturbation in the interface width arbitrarily small as this perturbation tends to zero. A similar argument applies in the discrete case. Even though the derivative is never actually zero, it suffices for this derivative to approach zero at a linear rate in order to be able to absorb a first-order error in the discrete interface width into a second-order error in the actual profile. \diamond

Remark 61. Regardless of the asymptotics themselves, an important (and quite natural) point to be retained from the previous remark is that the interface width by itself is an unreliable criterion for judging the actual quality of a solution. This is of course also (but therefore not only) an issue when dealing with well-potentials and might partially explain why it is quite common to “define” the interface width e.g. as being the width between the points where $\phi = 0.1$ and $\phi = 0.9$ instead of two points where ϕ is actually closer to 0 and 1. \diamond

Combining the estimates for κ and N , it further follows that

$$\kappa N = \frac{4}{\pi} \Delta\bar{x} \left(1 + \mathcal{O}((\Delta\bar{x})^2) \right) \left(\frac{\pi^2}{4\Delta\bar{x}} - (1 - \delta) + \mathcal{O}(\Delta\bar{x}) \right) = \pi - \frac{4\Delta\bar{x}}{\pi} (1 - \delta) + \mathcal{O}((\Delta\bar{x})^2),$$

and thus combined with $\cot\left(\frac{\pi}{2} - \theta\right) = \tan(\theta) = \theta + \mathcal{O}(\theta^3)$ that

$$\cot\left(\frac{\kappa N}{2}\right) = \cot\left(\frac{\pi}{2} - \left(\frac{2\Delta\bar{x}}{\pi}(1 - \delta) + \mathcal{O}((\Delta\bar{x})^2)\right)\right) = \frac{2\Delta\bar{x}}{\pi}(1 - \delta) + \mathcal{O}((\Delta\bar{x})^2). \quad (6.53)$$

³⁷E.g. if interpreted in a cell-centered fashion, there would be $N - 1$ “cells” with values deviating from the bulk-ones, whereas, if interpreted in a FEM-like fashion, there would be N such “elements”.

³⁸Multiplying κ by the “point” i leads to value

$$\kappa i = \frac{4}{\pi\epsilon} i \Delta x (1 + \mathcal{O}((\Delta x)^2)) = \frac{4}{\pi\epsilon} x_i (1 + \mathcal{O}((\Delta x)^2))$$

for the argument of the sine and cosine in Equation (6.36).

Remark 62. ◇

A multiplication of the second expression for $\mathcal{E}(N)$ in Equation (6.46) by Δx then shows that

$$\begin{aligned}\Delta x \mathcal{E} &= \frac{4\gamma}{\pi^2} N \Delta \bar{x} + \frac{4\gamma \Delta \bar{x}}{\pi^2} \cot\left(\frac{\kappa N}{2}\right) \sqrt{\frac{\pi^2}{4(\Delta \bar{x})^2} - 1} = \frac{4\gamma}{\pi^2} N \Delta \bar{x} + \frac{2\gamma}{\pi} \cot\left(\frac{\kappa N}{2}\right) \sqrt{1 - \frac{4(\Delta \bar{x})^2}{\pi^2}} \\ &= \frac{4\gamma}{\pi^2} \left(\frac{\pi^2}{4} - (1 - \delta) \Delta \bar{x} + \mathcal{O}((\Delta \bar{x})^2) \right) + \frac{2\gamma}{\pi} \left(\frac{2\Delta \bar{x}}{\pi} (1 - \delta) + \mathcal{O}((\Delta \bar{x})^2) \right) \left(1 - \mathcal{O}((\Delta \bar{x})^2) \right).\end{aligned}$$

From this, it is seen that the first-order errors through the $(1 - \delta)$ -terms actually cancel (up to second order) in the sum, leaving the final estimate

$$\Delta x \mathcal{E} = \gamma + \mathcal{O}((\Delta \bar{x})^2), \quad (6.54)$$

which in particular shows that the discrete energy is **second-order** convergent to the continuous value.

Remark 63. This very pleasant conclusion is in fact not obvious based on the use of a formally second-order numerical scheme. It is not a priori clear that this accuracy can actually be achieved due to the only moderate smoothness in relation with the jumps in the second derivatives at the transition points.

In continuation of Remark 51, one can also observe that, even though it is the first term in Equation (6.46) that actually converges to the desired value γ , it does so only with first-order accuracy due to the $\mathcal{O}(\frac{\Delta x}{\epsilon})$ -error in the discrete interface width $N\Delta x$. That the convergence of the energy is nevertheless quadratic is only achieved through the correction by the second term in Equation (6.46) which (up to second-order order) cancels this contribution.

It can also be noted that this does not depend on the precise form of the error in $N\Delta x$ (provided the error δN is of order 1, as e.g. the one caused by the truncation in the calculation of N). This can of course be seen by a close look at the coefficients arising in the expansions above, but is intuitively much easier to understand based upon the choice of N itself. In fact, the energetic argument underlying the choice of the “correct” number of points in the discrete interface is simply a discrete analogue of the continuous condition $\frac{\partial \mathcal{E}}{\partial N}(N) \stackrel{!}{=} 0$. Even though this equality need not be achieved exactly in the discrete setting, it will be so in first order. From this, the “direct” contribution $\Delta x \frac{\partial \mathcal{E}}{\partial N} \delta N$ of a perturbation in N to the discrete energy is indeed expected to be only of second-order in Δx . ◇

Remark 64. A further interesting observation is, by the expressions for the sums of the discrete bulk potential and gradient energy in equations (6.42) and (6.43), that there is - in contrast to the continuous case - not a precise equipartition of energy. In fact, one has

$$\sum_{i=1}^{N-1} \left(\epsilon a_i - \frac{1}{\epsilon} w_i \right) = \frac{4\gamma}{\pi^2 \epsilon} \left(\cot^2\left(\frac{\kappa N}{2}\right) N + \cot\left(\frac{\kappa N}{2}\right) \frac{1 - \cos(\kappa)}{\sin(\kappa)} \right),$$

showing that the gradient energy is larger than the bulk potential energy in the discrete case. Even though one could then conjecture that the second-order accuracy of $\Delta x \mathcal{E}$ is due to a beneficial cancellation, a short calculation using the expansions above reveals that both the gradient energy and bulk energy contributions themselves (and not just their sum $\Delta x \mathcal{E}$) are also second-order accurate³⁹. ◇

³⁹Both equations (6.42) and (6.43) contain a term of the form $\left(1 \pm \cot^2\left(\frac{\kappa N}{2}\right)\right)N$, which would be the obvious candidate for a cancellation. As $\cot\left(\frac{\kappa N}{2}\right) = \mathcal{O}(\Delta x)$ by the above, this term does not actually matter in the first order. More precisely, using this estimate for the cot and $\frac{1}{\sin(\kappa)} \approx \frac{1}{\kappa}$ (or directly inserting the definition of $\sin(\kappa)$)

As a final consequence of the expansions above, one can also construct a (quite accurate) continuous approximation to the discrete solution given in Equation (6.48) in terms of the parameters Δx and ϵ alone. In contrast to the exact discrete profile, this has the advantage of directly highlighting some simple - both qualitative and quantitative - points.

This requires additionally approximating the term $\frac{1}{\sin\left(\frac{\kappa N}{2}\right)}$. Based on equations (6.51) and (6.52), it follows (as already used above) that $\frac{\kappa N}{2} = \frac{\pi}{2} - \frac{2\Delta\bar{x}}{\pi}(1-\delta) + \mathcal{O}((\Delta\bar{x})^2)$. The elementary relation $\sin\left(\frac{\pi}{2} - \theta\right) = \cos(\theta) = 1 - \theta^2 + \mathcal{O}(\theta^4)$ near zero first leads to

$$\frac{1}{\sin\left(\frac{\kappa N}{2}\right)} = \frac{1}{1 - \frac{1}{2}\left(\frac{2\Delta\bar{x}}{\pi}(1-\delta)\right)^2 + \mathcal{O}((\Delta\bar{x})^3)} = 1 + \frac{1}{2}\left(\frac{2\Delta\bar{x}}{\pi}(1-\delta)\right)^2 + \mathcal{O}((\Delta\bar{x})^3)$$

and thus with $\kappa = \frac{4}{\pi\epsilon}(1 + \mathcal{O}((\Delta\bar{x})^2))\Delta x$ from Equation (6.51) shows that the discrete profile can also be rewritten as

$$\phi_i = \frac{1}{2} + \frac{1}{2}\left(1 + \frac{1}{2}\left(\frac{2\Delta\bar{x}}{\pi}(1-\delta)\right)^2\right) \sin\left(\frac{4}{\pi\epsilon}(1 + \mathcal{O}((\Delta\bar{x})^2))(x_i - x_m)\right) + \mathcal{O}((\Delta\bar{x})^3) \quad (6.55)$$

where $x_i := i\Delta x$.

Expanding the ‘‘error’’ in the angular frequency leads to

$$\begin{aligned} \sin\left(\frac{4}{\pi\epsilon}(1 + \mathcal{O}((\Delta\bar{x})^2))(x_i - x_m)\right) &= \sin\left(\frac{4}{\pi\epsilon}(x_i - x_m)\right) \cos\left(\mathcal{O}((\Delta\bar{x})^2)(x_i - x_m)\right) \\ &\quad + \cos\left(\frac{4}{\pi\epsilon}(x_i - x_m)\right) \sin\left(\mathcal{O}((\Delta\bar{x})^2)(x_i - x_m)\right). \end{aligned} \quad (6.56)$$

Since $\cos\left(\mathcal{O}((\Delta\bar{x})^2)\right) = 1 - \mathcal{O}((\Delta\bar{x})^4)$, the prefactor in the sine-term can actually be replaced (up to a fourth-order error) by 1, whereas the linear convergence of $\sin(\theta)$ to zero with θ in principle allows dropping the second-order term up to an error of $\mathcal{O}((\Delta x)^2)$. Together with the boundedness of the sin- and cos-function, this is already sufficient to show that the discrete profile satisfies

$$\phi_i = \frac{1}{2} + \frac{1}{2} \sin\left(\frac{4}{\pi\epsilon}(x_i - x_m)\right) + \mathcal{O}((\Delta\bar{x})^2)$$

and is thus indeed a second-order accurate approximation of the analytical one.

Remark 65. It should again be stressed that this is not a priori obvious just because one is using a formally second-order scheme for evaluating $\frac{d^2\phi}{dx^2}$. Even though imposing e.g. the requirement

in Equation (6.33), one has

$$\begin{aligned} \Delta x \sum_{i=1}^{N-1} \epsilon a_i &= \frac{8\gamma}{\pi^2} \Delta\bar{x} \left(\frac{1}{4} \left(1 + \cot^2\left(\frac{\kappa N}{2}\right)\right) N + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \frac{1}{\sin(\kappa)} \right) = \frac{8\gamma}{\pi^2} \Delta\bar{x} \left(\frac{1}{4} \left(1 + \mathcal{O}((\Delta\bar{x})^2)\right) N + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \frac{1}{\sin(\kappa)} \right) \\ &= \frac{8\gamma}{\pi^2} \Delta\bar{x} \left(\frac{1}{4} \left(1 + \mathcal{O}((\Delta\bar{x})^2)\right) \left(\frac{\pi^2}{4\Delta\bar{x}} - (1-\delta) + \mathcal{O}(\Delta\bar{x}) \right) + \frac{1}{2} \left(\frac{2\Delta\bar{x}}{\pi} (1-\delta) + \mathcal{O}((\Delta\bar{x})^2) \right) \frac{\pi}{4\Delta\bar{x}} \left(1 + \mathcal{O}((\Delta\bar{x})^2)\right) \right) \\ &= \frac{8\gamma}{\pi^2} \Delta\bar{x} \left(\frac{\pi^2}{16\Delta\bar{x}} - \frac{1}{4}(1-\delta) + \frac{1}{4}(1-\delta) + \mathcal{O}(\Delta\bar{x}) \right) = \frac{1}{2} \gamma + \mathcal{O}((\Delta\bar{x})^2). \end{aligned}$$

The contribution by w differs, up to the sign of the second-order term $\cot^2\left(\frac{\kappa N}{2}\right)$ only by the last factor $\cot(\kappa)$ instead of $\frac{1}{\sin(\kappa)}$. As $\cot(\kappa) = \frac{\cos(\kappa)}{\sin(\kappa)} = \frac{1}{\kappa} + \mathcal{O}((\Delta\bar{x})^2)$, the calculation for the contributions by a can essentially directly be reused to show that one also has

$$\Delta x \sum_{i=1}^{N-1} \frac{1}{\epsilon} w(\phi_i) = \frac{1}{2} \gamma + \mathcal{O}((\Delta\bar{x})^2).$$

An even quicker alternative is of course to observe that as the gradient energy contribution and the total energy are second order accurate, the bulk potential necessarily also has to be as the difference between both.

$\phi \in C^4$ for obtaining the formal estimate $\frac{\phi_{i+1}-2\phi_i+\phi_{i-1}}{(\Delta x)^2} = \phi''(x_i) + \mathcal{O}((\Delta x)^2) + \mathcal{O}((\Delta x)^3)$ is in fact more restrictive than actually necessary for obtaining global second-order convergence, it should be kept in mind that, as pointed out in Remark 48, that this difference operator is only $\mathcal{O}(1)$ -accurate when it straddles the transition points, corresponding to a loss of two order of accuracy. Obtaining estimates for the order of convergence under such conditions requires a much more detailed and complex analysis (see e.g. [38] for a detailed treatment of this topic). \diamond

While this is interesting in itself, it is also useful to have a more precise estimate of the remaining second-order deviations from the analytical profile. As the analysis above shows, the accuracy of the approximation above is actually of the order $\mathcal{O}((\Delta \bar{x})^3)$ except for the ‘‘imprecise’’ second-order estimate for $\sin\left(\mathcal{O}((\Delta \bar{x})^2)(x_i - x_m)\right)$ due to the only second-order approximation of κ .

Improving this estimate is fortunately only a matter of solving a quadratic equation. In fact, κ being defined by $\cos(\kappa) = 1 - \frac{1}{2}\left(\frac{4\Delta \bar{x}}{\pi \epsilon}\right)^2$, one can simply include the fourth-order term in the expansion of the cosine to obtain $1 - \frac{1}{2}\kappa^2 + \frac{1}{4!}\kappa^4 + \mathcal{O}(\kappa^6) \stackrel{!}{=} 1 - \frac{1}{2}\left(\frac{4\Delta \bar{x}}{\pi}\right)^2$ resp.

$$\kappa^4 - 12\kappa^2 + 12\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}(\kappa^6) \stackrel{!}{=} 0.$$

Since $\kappa = \mathcal{O}(\Delta \bar{x})$, this can be interpreted as a quadratic equation for κ^2 perturbed by a term of order $(\Delta \bar{x})^6$, from which it follows (using $\sqrt{1+y} = 1 + \frac{1}{2}y - \frac{1}{8}y^2 + \mathcal{O}(y^3)$ near $y = 0$ and choosing the relevant sign preceding the square-root) that

$$\begin{aligned} \kappa^2 &= 6 - 6\sqrt{1 - \frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)} \\ &= 6 - 6\left(1 - \frac{1}{2}\left(\frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)\right) - \frac{1}{8}\left(\frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)\right)^2\right) \\ &= 6\left(\frac{1}{2}\left(\frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)\right) + \frac{1}{8}\left(\frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)\right)^2\right) \\ &= \left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6) + \frac{3}{4}\left(\frac{1}{3}\left(\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^6)\right)\right)^2 = \left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \frac{1}{12}\left(\frac{4\Delta \bar{x}}{\pi}\right)^4 + \mathcal{O}((\Delta \bar{x})^6). \end{aligned}$$

Again expanding the square-root, a more precise estimate for κ is therefore given by

$$\begin{aligned} \kappa &= \frac{4\Delta \bar{x}}{\pi} \sqrt{1 + \frac{1}{12}\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^4)} = \frac{4\Delta \bar{x}}{\pi} \left(1 + \frac{1}{24}\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^4)\right) \\ &= \frac{4\Delta \bar{x}}{\pi} + \frac{\Delta \bar{x}}{6\pi} \left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^5), \end{aligned} \tag{6.57}$$

allowing to obtain the improved approximation

$$\begin{aligned} &\cos\left(\frac{4}{\pi \epsilon}(x_i - x_m)\right) \sin\left(\left(\frac{1}{6\pi \epsilon}\left(\frac{4\Delta \bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta \bar{x})^4)\right)(x_i - x_m)\right) \\ &= \frac{1}{6\pi \epsilon} \left(\frac{4\Delta \bar{x}}{\pi}\right)^2 (x_i - x_m) \cos\left(\frac{4}{\pi \epsilon}(x_i - x_m)\right) + \mathcal{O}((\Delta \bar{x})^4) \end{aligned}$$

for the cos-based term in Equation (6.56). Combining this with the previous analysis for Equation

(6.55) leads to the higher-order approximation⁴⁰

$$\begin{aligned}\phi_i = & \frac{1}{2} + \frac{1}{2} \left(1 + \frac{1}{2} \left(\frac{2\Delta\bar{x}}{\pi} (1-\delta) \right)^2 \right) \sin \left(\frac{4}{\pi\epsilon} (x_i - x_m) \right) \\ & + \frac{1}{12\pi\epsilon} \left(\frac{4\Delta\bar{x}}{\pi} \right)^2 (x_i - x_m) \cos \left(\frac{4}{\pi\epsilon} (x_i - x_m) \right) + \mathcal{O}(\Delta\bar{x})^3.\end{aligned}\tag{6.58}$$

In contrast to the previous approximation in Equation (6.55), this expansion contains both leading-order (i.e. second-order) perturbations of the analytical profile.

The first one in terms of the slightly increased amplitude of the dominant sine-term was already obtained before. It on the one hand leads to a slightly steeper profile in the region around the isoline $\phi = \frac{1}{2}$, and on the other hand to both a slightly increased (resp. decreased) ϕ -value and a slightly decreased (resp. increased) curvature near the outer limits of the interface near $\phi = 1$ (resp. $\phi = 0$) since this is where both the sin itself and thus also its second derivative attain the maximal magnitude. The second contribution in contrast will also lead to a steeper profile near the middle of the interface (the cosine being second-order close to one there), but has very little effect on the ϕ -values near the outer limits as this is where the cosine approaches zero, its primary effect on the ϕ -values itself therefore lying in the intermediate region. Nevertheless, with $(x \cos(x))'' = -2 \sin(x) - x \cos(x)$, it will also have a similiar ‘‘curvature-enhancing’’ effect near the outer regions.

Remark 66. With respect to the interpretation of both perturbations and their relation with the **first-order** perturbation of the interface width, it is worth noting that the first contribution could also be obtained by simply enforcing a steeper version (for simplicity with $x_m = 0$) of the analytical profile $\frac{1}{2} + \frac{1}{2}c' \sin\left(\frac{\pi}{4\epsilon}x\right)$ to match the bulk-value 1 at the perturbed point $x' = \frac{\pi^2}{8}\epsilon - \frac{1-\delta}{2}\Delta\bar{x}$, corresponding to distributing the discrepancy in the width symmetrically to both sides. In fact, with

$$\sin\left(\frac{4}{\pi\epsilon}x'\right) = \sin\left(\frac{\pi}{2} - \frac{4}{\pi\epsilon}\frac{1-\delta}{2}\Delta\bar{x}\right) = \cos\left(\frac{2}{\pi\epsilon}(1-\delta)\Delta\bar{x}\right),$$

the equation $\frac{1}{2} + c' \sin\left(\frac{4}{\pi\epsilon}x'\right) \stackrel{!}{=} 1$ asymptotically reduces to

$$c' \cos\left(\frac{2}{\pi\epsilon}(1-\delta)\Delta\bar{x}\right) = c' \left(1 - \frac{1}{2} \left(\frac{2}{\pi\epsilon}(1-\delta)\Delta\bar{x} \right)^2 \right) + \mathcal{O}(\Delta\bar{x})^4 \stackrel{!}{=} \frac{1}{2}$$

and thus

$$c' = \frac{1}{2 \left(1 - \frac{1}{2} \left(\frac{2}{\pi\epsilon}(1-\delta)\Delta\bar{x} \right)^2 \right) + \mathcal{O}(\Delta\bar{x})^4} = \frac{1}{2} \left(1 + \frac{1}{2} \left(\frac{2}{\pi\epsilon}(1-\delta)\Delta\bar{x} \right)^2 \right) + \mathcal{O}(\Delta\bar{x})^4.$$

This is, up to a fourth-order error, the prefactor of the sin above. Note that, as already indicated before, it is crucial here that the first-order perturbation in the interface width arises in a region where the profile is to first order flat as this is the reason c' differs from 1 only in second and not in first order.

In contrast, the second contribution can be interpreted as being related to a more ‘‘mundane’’ discretization error in line with Footnote 28. More precisely, applying a discrete (three-point

⁴⁰Even though the perturbation by the term $(x_i - x_m) \cos\left(\frac{4}{\pi\epsilon}(x_i - x_m)\right)$ a priori also appears with the slightly increased amplitude of the sin-term, this results only in a fourth-order perturbation.

stencil) second-order derivative to a sinusoidal profile, one obtains

$$\begin{aligned}
& \frac{\sin(x_{i+1}) - 2\sin(x_i) + \sin(x_{i-1}))}{(\Delta x)^2} = \frac{\sin(x_i + \Delta x) - 2\sin(x_i) + \sin(x_i - \Delta x)}{(\Delta x)^2} \\
& = \frac{\sin(x_i) \cos(\Delta x) + \cos(x_i) \sin(\Delta x) - 2\sin(x_i) + \sin(x_i) \cos(\Delta x) - \cos(x_i) \sin(\Delta x)}{(\Delta x)^2} \\
& = -\sin(x_i) \frac{2(1 - \cos(\Delta x))}{(\Delta x)^2} = -\sin(x_i) \frac{2\left(1 - \left(1 - \frac{1}{2}(\Delta x)^2 + \frac{1}{4!}(\Delta x)^4 - \mathcal{O}((\Delta x)^6)\right)\right)}{(\Delta x)^2} \\
& = -\left(1 - \frac{1}{12}(\Delta x)^2\right) \sin(x_i) + \mathcal{O}((\Delta x)^4),
\end{aligned}$$

resp. from the same calculation but with a modified angular frequency for the sine,

$$\frac{\sin\left(\frac{4}{\pi\epsilon}x_{i+1}\right) - 2\sin\left(\frac{4}{\pi\epsilon}x_i\right) + \sin\left(\frac{4}{\pi\epsilon}x_{i-1}\right)}{(\Delta x)^2} = -\left(\frac{4}{\pi\epsilon}\right)^2 \sin\left(\frac{4}{\pi\epsilon}x_i\right) \left(1 - \frac{1}{12}\left(\frac{4\Delta\bar{x}}{\pi}\right)^2 + \mathcal{O}((\Delta\bar{x})^4)\right).$$

One can successfully counterbalance this deviation in the second derivative of the dominant term precisely by the cos-based term in Equation (6.58). In fact, inserting the ansatz (again taking $x_m = 0$ for simplicity) $\phi(x) = \frac{1}{2} + \frac{1}{2}\sin\left(\frac{4}{\pi\epsilon}x\right) + \delta\phi$ with $\delta\phi = \mathcal{O}((\Delta\bar{x})^2)$ into the discrete difference Equation (6.31), it follows that $\delta\phi$ should satisfy

$$\begin{aligned}
& -2\gamma\epsilon \left(-\left(\frac{4}{\pi\epsilon}\right)^2 \frac{1}{2} \left(1 - \frac{1}{12}\left(\frac{4\Delta\bar{x}}{\pi}\right)^2\right) \sin\left(\frac{4}{\pi\epsilon}x\right) - 2\gamma\epsilon \frac{d^2\delta\phi}{dx^2} + \frac{16}{\pi^2\epsilon} \gamma \left(-\sin\left(\frac{4}{\pi\epsilon}x\right) - 2\delta\phi\right) \right. \\
& \left. = -\gamma\epsilon \left(\left(\frac{4}{\pi\epsilon}\right)^2 \frac{1}{12} \left(\frac{4\Delta\bar{x}}{\pi}\right)^2 \sin\left(\frac{4}{\pi\epsilon}x\right) - 2\gamma\epsilon \frac{d^2\delta\phi}{dx^2} - \frac{32}{\pi^2\epsilon} \gamma \delta\phi \right) \approx 0
\end{aligned}$$

since the first-order terms in the sin cancel by construction of the analytical profile and applying the discrete second-order derivative to $\delta\phi$ will approximate its actual second derivative in order $\mathcal{O}((\Delta\bar{x})^2)$. $\delta\phi$ should thus satisfy

$$-2\epsilon \frac{d^2\delta\phi}{dx^2} - \frac{32}{\pi^2\epsilon} \delta\phi \approx \frac{4}{3\pi^2\epsilon} \left(\frac{4\Delta\bar{x}}{\pi}\right)^2 \sin\left(\frac{4}{\pi\epsilon}x\right)$$

and it is easily seen that defining $\delta\phi := \frac{1}{12\pi\epsilon} \left(\frac{4\Delta\bar{x}}{\pi}\right)^2 x \cos\left(\frac{4}{\pi\epsilon}x\right)$ corresponding to the second correction in Equation (6.58) will actually solve this equation⁴¹.

It can further be noted that, for the choice of κ given by Equation (6.33), one has

$$\begin{aligned}
& \frac{\sin(\kappa(i+1)) - 2\sin(\kappa i) + \sin(\kappa(i-1)))}{(\Delta x)^2} = -2\sin(\kappa i) \frac{1 - \cos(\kappa)}{(\Delta x)^2} \\
& = -2\sin(\kappa i) \frac{1 - \left(1 - \frac{1}{2}\left(\frac{4\Delta x}{\pi\epsilon}\right)^2\right)}{(\Delta x)^2} = -\left(\frac{4}{\pi\epsilon}\right)^2 \sin(\kappa i),
\end{aligned}$$

i.e. the prefactor of the sine resulting from the discrete second-order differentiation operator coincides **exactly** with the one arising from the second differentiation of the $\sin\left(\frac{4}{\pi\epsilon}x\right)$ in the continuous solution. \diamond

⁴¹Note that while this $\delta\phi$ solves the equation exactly in the interior of the domain, it is not completely compatible with the (here implicit) boundary condition that $\delta\phi$ should vanish at the endpoints of the actual interface. This violation is only first-order in Δx (with respect to the already second-order size of $\delta\phi$ itself) though by the linear convergence of the cos-term to zero near the outer region. This can then either be absorbed through the construction of an appropriate boundary layer, or, since the dominant sin-term is maximal there, by a slight modification of the dominant term through a third-order correction of its amplitude. Regardless of the details, this does not affect the order of the error estimate in Equation (6.58).

6.3 Multi-Phasefield Problems

As discussed in the beginning of Section 6.2, one has to deal with a number of additional challenges when considering problems with more than two phases.

A first - and quite important one within the usual variational framework - is the construction of an appropriate generalization of the underlying functional to the multiphase case. While there are a variety of possible choices, the discussion here will focus on the one based on [52] outlined in Section 6.1. Together with the additional constraints through either the sum-constraint $\sum_{\alpha=1}^N \phi^\alpha = 1$ or the restriction of the ϕ -values to the Gibbs-simplex \mathcal{GS}^N , the functional already fixes an important part of the problem, namely the equations to be satisfied by any minimizer and thus the steady-state equations, which will quickly be rederived - in a formal manner similar to Section 6.2 - in Subsection 6.3.1.

A second important aspect, which is essentially independent of the functional (even if one postulates a gradient flow)⁴² is the choice of the dynamics, i.e. how one chooses the “proportionality” with respect to $\frac{\delta \mathcal{F}_\epsilon}{\delta \phi}$. Two popular choices also implemented within the **Pace3D**-framework will quickly be outlined in Subsection 6.3.2. Despite its practical importance, the (relatively difficult) question of the respective advantages and disadvantages of the different choices of dynamics from a physical point of view will not be discussed in any detail here. Instead, after a short outline of some general algorithmic considerations for multiphase problems in Subsection 6.3.3, in particular in the presence of a large number of phases, the focus in Subsection 6.3.3 will primarily be on practical aspects of this choice in the presence of the constraint by the Gibbs-simplex.

6.3.1 The Steady-State Equations

As in the two-phase case, a directional differentiation of the phasefield-functional in Equation (6.10) leads to the first-order necessary condition

$$\mathcal{F}'_\epsilon(\phi; \psi) = \sum_{\alpha=1}^N \int_{\Omega} \epsilon \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \nabla \psi^\alpha + \frac{\partial a}{\partial \phi^\alpha} \psi^\alpha \right) + \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} \psi^\alpha + \frac{\partial f}{\partial \phi^\alpha} \psi^\alpha \, d\mathbf{x} \geq 0$$

for all admissible directions ψ .

If the individual contributions to \mathcal{E}_ϵ and the functions themselves are smooth enough and with the “natural” boundary condition $\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \mathbf{n} = 0$ (see Remark 70 below), any minimizer in the multi-well case (without an additional restriction of the form $0 \leq \phi^\alpha \leq 1$) will therefore have to satisfy

$$\sum_{\alpha=1}^N \int_{\Omega} \underbrace{\left(-\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) + \epsilon \frac{\partial a}{\partial \phi^\alpha} + \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} + \frac{\partial f}{\partial \phi^\alpha} \right)}_{=: g^\alpha} \psi^\alpha \, d\mathbf{x} \geq 0, \quad (6.59)$$

where, unlike in Equation (6.17), ψ is locally restricted to satisfy $\sum_{\alpha=1}^N \psi^\alpha = 0$. In the steady-state, the (L^2) -gradients $\mathbf{g} = (g^\alpha)_{1 \leq \alpha \leq N}$ are therefore locally orthogonal to all vectors with zero average, i.e., as seen in Section 4, satisfy $\mathbf{g}(\mathbf{x}) = -\Lambda(\mathbf{x})\mathbf{e}$, where the (scalar) Lagrange-multiplier $\Lambda(\mathbf{x})$ does depend on \mathbf{x} but not on the phase α and the choice of the minus-sign is only a matter of convenience. Through a simple summation, it therefore follows that Λ satisfies $\sum_{\alpha=1}^N -g^\alpha = N\Lambda$, and thus, inserting the expression for \mathbf{g} , that ϕ has to satisfy as in [52]

$$\begin{cases} -\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) + \epsilon \frac{\partial a}{\partial \phi^\alpha} + \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} + \frac{\partial f}{\partial \phi^\alpha} + \Lambda - \mu^\alpha = 0, \\ \sum_{\alpha=1}^N \phi^\alpha = 1 \end{cases} \quad (6.60)$$

⁴²It is partially related to the constraint sets though.

together with the appropriate isolating boundary conditions and Λ locally given by

$$\Lambda = \frac{1}{N} \sum_{\alpha=1}^N \epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) - \epsilon \frac{\partial a}{\partial \phi^\alpha} - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} - \frac{\partial f}{\partial \phi^\alpha}. \quad (6.61)$$

If the phasefield values are additionally subject to the constraints $0 \leq \phi^\alpha \leq 1 \forall \alpha$, the conditions are, due to the additional sum-constraint, somewhat more complex than in the two-phase setting from Section 6.2. Considering first the simpler representation of the (locally) admissible set in Equation (6.2), one formally has that ψ is an admissible local variation if ψ satisfies $\sum_{\alpha=1}^N \psi^\alpha = 0$ with the additional restriction that $\psi^\alpha \geq 0$ if $\phi^\alpha = 0$. As in e.g. [11], if there are at least two **free phases** with $0 < \phi^\alpha, \phi^\beta < 1$, taking any $\psi > 0$ and $\psi = \pm \psi e^\alpha \mp \psi e^\beta$ - thus satisfying the sum-constraint - Equation (6.59) shows that $\pm \psi (g^\alpha - g^\beta) \geq 0$ and therefore $g^\alpha = g^\beta$. As this argument can be repeated for an arbitrary combination of free phases, it follows that

$$g^\alpha(\mathbf{x}) = -\Lambda(\mathbf{x}) \quad \forall \alpha : 0 < \phi^\alpha < 1, \quad (6.62)$$

where Λ is, just as in the simpler case when restricted to Σ_1^N only, independent of the free phases. If there is any phase α at 0, there is necessarily at least one phase β with $\phi^\beta > 0$ due to the sum-constraint. In this case, one has to restrict the sign of the variations of ϕ^α taking e.g. again $\psi > 0$ and $\psi = \psi e^\alpha - \psi e^\beta$, and one can only conclude that $\psi (g^\alpha - g^\beta) \geq 0$, i.e. $g^\alpha \geq g^\beta$ if $\phi^\alpha = 0$, $\phi^\beta > 0$.

If there is any β such that $0 < \phi^\beta < 1$, by the condition above, one therefore has $g^\alpha \geq -\Lambda$. It may happen though that the (then necessarily only) phase with $\phi^\beta > 0$ actually satisfies $\phi^\beta = 1$, in which case all other phases are automatically at 0. As there are no free phases, the above definition of Λ does not apply. Instead, one only has $g^\alpha \geq g^\beta \forall \alpha \neq \beta$. If one in this case **chooses** $\Lambda := -g^\beta$, the necessary conditions can be written in a unified manner in terms of the KKT system

$$\mathbf{g} = -\Lambda \mathbf{e} + \boldsymbol{\mu}, \quad \boldsymbol{\mu} \geq 0, \quad \mu^\alpha \phi^\alpha = 0, \quad (6.63)$$

or, more explicitly,

$$\begin{cases} -\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) + \epsilon \frac{\partial a}{\partial \phi^\alpha} + \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} + \frac{\partial f}{\partial \phi^\alpha} + \Lambda - \mu^\alpha = 0, \\ \sum_{\alpha=1}^N \phi^\alpha = 1, \\ \mu^\alpha \geq 0, \\ \mu^\alpha \phi^\alpha = 0 \end{cases} \quad (6.64)$$

for $\alpha = 1, \dots, N$.

In contrast to the simpler well-case, there is now **no** simple explicit formula for Λ as in Equation (6.61) anymore. In fact, the same summation procedure as in the well-case applied to the first line in Equation (6.63) now shows that

$$\Lambda = \underbrace{\frac{1}{N} \sum_{\alpha=1}^N \left(\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) - \epsilon \frac{\partial a}{\partial \phi^\alpha} - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} - \frac{\partial f}{\partial \phi^\alpha} \right)}_{=: \Lambda_{mw}} + \frac{1}{N} \sum_{\alpha=1}^N \mu^\alpha = \Lambda_{mw} + \frac{1}{N} \sum_{\alpha=1}^N \mu^\alpha. \quad (6.65)$$

Since the μ^α are non-negative, the expression for Λ in Equation (6.61) is therefore only a lower bound for the actual value of Λ in the obstacle case.

Remark 67. As the a -function in [52] is generally only defined indirectly in terms of ϕ and $\nabla \phi$ based upon the generalized gradient vectors $\mathbf{q}^{\alpha\beta}(\phi, \nabla \phi) = \phi^\alpha \nabla \phi^\beta - \phi^\beta \nabla \phi^\alpha$ from Equation (6.5), it is often more convenient to instead consider a as a function of the $\mathbf{q}^{\alpha\beta}$ only and then obtain the derivative with respect to ϕ^α and $\nabla \phi^\alpha$ through a simple chain rule. Since⁴³

$$\frac{\partial \mathbf{q}^{\beta\eta}}{\partial \phi^\alpha} = \nabla \phi^\eta \delta^{\alpha\beta} - \nabla \phi^\beta \delta^{\alpha\eta} \quad \text{and} \quad \frac{\partial \mathbf{q}^{\beta\eta}}{\partial \nabla \phi^\alpha} = (\phi^\beta \delta^{\alpha\eta} - \phi^\eta \delta^{\alpha\beta}) \mathbf{I},$$

⁴³Where $\delta^{\alpha\beta}$ is the usual Kronecker symbol, i.e. 1 if $\alpha = \beta$ and 0 otherwise.

the differentiation of the summation in Equation (6.6) leads to

$$\frac{\partial a}{\partial \phi^\alpha} = \frac{\partial}{\partial \phi^\alpha} \sum_{\beta} \sum_{\eta > \beta} A^{\beta\eta}(\mathbf{q}^{\beta\eta}) = \sum_{\beta} \sum_{\eta > \beta} \frac{\partial A^{\beta\eta}}{\partial \mathbf{q}^{\beta\eta}} (\nabla \phi^\eta \delta^{\alpha\beta} - \nabla \phi^\beta \delta^{\alpha\eta}).$$

The summation over the first term then directly reduces to $\sum_{\eta > \alpha} \frac{\partial A^{\alpha\eta}}{\partial \mathbf{q}^{\alpha\eta}} \cdot \nabla \phi^\eta$, whereas the second one does, after exchanging the order of summation through $\sum_{\beta} \sum_{\eta > \beta} (\cdot)^{\beta\eta} = \sum_{\eta} \sum_{\beta < \eta} (\cdot)^{\beta\eta}$, lead to $-\sum_{\beta < \alpha} \frac{\partial A^{\beta\alpha}}{\partial \mathbf{q}^{\beta\alpha}} \cdot \nabla \phi^\beta$. With $\mathbf{q}^{\alpha\beta} = -\mathbf{q}^{\beta\alpha}$ and $A^{\alpha\beta}(\mathbf{q}^{\alpha\beta}) = A^{\beta\alpha}(\mathbf{q}^{\beta\alpha}) = A^{\beta\alpha}(-\mathbf{q}^{\alpha\beta})$ and exchanging dummy-indices, one obtains the “missing” other half of the summation and thus finally

$$\frac{\partial a(\phi, \nabla \phi)}{\partial \phi^\alpha} = \sum_{\beta \neq \alpha} \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}} \cdot \nabla \phi^\beta. \quad (6.66)$$

The same argument applied to the differentiation with respect to $\nabla \phi^\alpha$ in turn shows that

$$\frac{\partial a(\phi, \nabla \phi)}{\partial \nabla \phi^\alpha} = - \sum_{\beta \neq \alpha} \phi^\beta \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}. \quad (6.67)$$

This in principle simple but still somewhat tedious reasoning occurs very frequently in the multiphase case and will for the sake of brevity not be repeated in detail again⁴⁴, but can be verified to lead to the expressions

$$\frac{\partial w_{mw}(\phi)}{\partial \phi^\alpha} = 18\phi^\alpha \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} (\phi^\beta)^2 + 2\phi^\alpha \sum_{\substack{\delta > \beta \\ \beta, \delta \neq \alpha}} \gamma^{\alpha\beta\delta} (\phi^\beta)^2 (\phi^\delta)^2 \quad (6.68)$$

and

$$\frac{\partial w_{mo}(\phi)}{\partial \phi^\alpha} = \frac{16}{\pi^2} \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \phi^\beta + \sum_{\substack{\delta > \beta \\ \beta, \delta \neq \alpha}} \gamma^{\alpha\beta\delta} \phi^\beta \phi^\delta \quad (6.69)$$

for the derivative of the bulk-potentials in the multi-well and multi-obstacle case. \diamond

Remark 68. The choice $\Lambda = -g^\beta$ for the bulk-phase β above is the most natural one when explicitly enforcing only the positivity constraint explicitly and the constraint $\phi^\alpha \leq 1$, $\alpha = 1, \dots, N$ implicitly through the sum-constraint as in the description of the Gibbs-simplex in Equation (6.2). One could of course, as in Equation (6.1), also enforce this upper bound explicitly, thus introducing a multiphase equivalent of the two multipliers μ^\pm in Section 6.2, then leading to the multiphase analogon

$$\begin{cases} -\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) + \epsilon \frac{\partial a}{\partial \phi^\alpha} + \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} + \frac{\partial f}{\partial \phi^\alpha} + \Lambda + \mu^{+, \alpha} - \mu^{-, \alpha} = 0, \\ \mu^{\pm, \alpha} \geq 0, \\ \mu^{-, \alpha} \phi^\alpha = 0, \\ \mu^{+, \alpha} (1 - \phi^\alpha) = 0 \end{cases}$$

of Equation (6.23).

This difference is of course only relevant when there is actually a bulk phase, but then has the disadvantage of introducing an indeterminacy due to the partial redundancy of the description of the constraint set. In fact, the sign restriction on the multipliers μ^\pm in the presence of a bulk-phase could be satisfied for any Λ satisfying $g^\beta \leq -\Lambda \leq \min_{\alpha \neq \beta} g^\alpha$ as for any such value, one would still have $g^\alpha = -\Lambda + \mu^{-, \alpha}$ with $\mu^{-, \alpha} \geq 0$ for all phases with $\phi^\alpha = 0$, whereas $g^\beta = -\Lambda - \mu^{+, \beta}$

⁴⁴A slightly quicker approach is to make use of the “symmetry” of the energy with respect to the pairings and thus to replace the summation over all distinct phase-pairings with one-half the summation over all phase-pairings.

with $\mu^{+, \beta} \geq 0$. Even though this indeterminacy is not a problem in principle and the particular choice of Λ (and thus $\mu^{+, \beta}$) does not affect the actual result, this is an unnecessary complication which can easily be avoided by using the non-redundant description in Equation (6.2).

It should be noted though that this leads to a difference of the multipliers when using a multiphase formulation with two phases as compared to the reduced two-phase formulation. The reduced formulation implicitly corresponds to choosing (up to a factor of 2) the multiplier $\Lambda = -\frac{1}{2}(g^1 + g^2)$ and thus will lead to the multipliers μ^\pm being equal to (twice) that value instead of the choice g^β above anytime a bound-constraint is active (which then automatically implies the presence of a bulk-phase)⁴⁵. \diamond

Remark 69. Similar indeterminacies arise quite frequently in non-reduced formulations. Another simple example in relation with a volume-constraint will be seen in Section 6.3.3, whereas a more complex one in a concentration-based setting will be discussed in some detail in Section 7.1. \diamond

Remark 70. Recall from Section 6.2 that imposing $\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \mathbf{n}$ is much less obvious in the presence of constraints than it might appear at first sight, since the boundary condition it is not a priori “up for choice” but actually has to be taken in accordance with the first-order necessary condition.

In particular, due to the multiphase setting, there is now even in the well-case an additional restriction $\sum_{\alpha=1}^N \psi^\alpha = 0$ on the variations ψ^α , which a priori only leads to the necessity of $\sum_{\alpha=1}^N \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \mathbf{n} \right) \psi^\alpha \stackrel{!}{=} 0$ and thus $\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \mathbf{n} = \Lambda$ for all α .

In contrast to the interior of the domain, where “choosing” to enforce $g^\alpha = 0$, $\alpha = 1, \dots, N$ (and thus $\Lambda = 0$) is potentially possible, but practically useless for making the volume integral vanish⁴⁶ (and therefore), enforcing $\Lambda = 0$ and therefore the vanishing of the conormal derivative on the boundary would seem a lot less restrictive. Nevertheless, even though this is clearly compatible with the first-order necessary condition, it is not a priori obvious whether this is a necessity or simply a selection having - with the pure bulk solution with the lowest f^α - at least one solution which is also a global minimizer. Showing that this is indeed a necessity is expected, but, similar to the discussion in Section 6.2, actually requires a non-trivial argument even in the simplest isotropic case⁴⁷ and will therefore simply be assumed here. \diamond

⁴⁵Note also that $g^\alpha \geq g^\beta$ in the two-phase case automatically ensures $g^\beta \leq -\Lambda = \frac{g^\alpha + g^\beta}{2} \leq g^\alpha$. A similar inequality based on the average of more than two phases however is generally **not** valid.

⁴⁶This is basically only compatible with the solution consisting out of a single bulk-phase either in the absence of driving forces or using one of the h -functions satisfying $h'(0) = h'(1) = 0$.

⁴⁷Based on simplifying the expression for (6.67), one has

$$\frac{\partial a}{\partial \nabla \phi^\alpha} \cdot \mathbf{n} = - \sum_{\beta \neq \alpha} \phi^\beta \mathbf{q}^{\alpha, \beta} \cdot \mathbf{n} = - \sum_{\beta \neq \alpha} \phi^\beta \left(\phi^\alpha \frac{\partial \phi^\beta}{\partial \mathbf{n}} - \phi^\beta \frac{\partial \phi^\alpha}{\partial \mathbf{n}} \right) \stackrel{!}{=} \Lambda, \alpha = 1, \dots, N.$$

Similar to the interior of the domain, a simple summation over all phases, but now using $\sum_{\beta \neq \alpha} = 1 - \phi^\alpha$ as well as $\sum_{\beta \neq \alpha} \frac{\partial \phi^\beta}{\partial \mathbf{n}} = -\frac{\partial \phi^\alpha}{\partial \mathbf{n}}$ due to the sum-constraint, shows that

$$\sum_{\beta} \sum_{\alpha \neq \beta} \phi^\beta \left(\phi^\alpha \frac{\partial \phi^\beta}{\partial \mathbf{n}} - \phi^\beta \frac{\partial \phi^\alpha}{\partial \mathbf{n}} \right) = \sum_{\beta} \phi^\beta \left((1 - \phi^\beta) \frac{\partial \phi^\beta}{\partial \mathbf{n}} - \phi^\beta \left(-\frac{\partial \phi^\beta}{\partial \mathbf{n}} \right) \right) = \sum_{\beta} \phi^\beta \frac{\partial \phi^\beta}{\partial \mathbf{n}} = -N\Lambda. \quad (6.70)$$

Making use of this relation, one can eliminate one of the summations in the condition on the normal derivative for the individual phases, since

$$\begin{aligned} & - \sum_{\beta \neq \alpha} \phi^\beta \left(\phi^\alpha \frac{\partial \phi^\beta}{\partial \mathbf{n}} - \phi^\beta \frac{\partial \phi^\alpha}{\partial \mathbf{n}} \right) = -\phi^\alpha \sum_{\beta \neq \alpha} \phi^\beta \frac{\partial \phi^\beta}{\partial \mathbf{n}} + \sum_{\beta \neq \alpha} (\phi^\beta)^2 \frac{\partial \phi^\alpha}{\partial \mathbf{n}} \\ & = -\phi^\alpha \left(\left(\sum_{\beta} \phi^\beta \frac{\partial \phi^\beta}{\partial \mathbf{n}} \right) - \phi^\alpha \frac{\partial \phi^\alpha}{\partial \mathbf{n}} \right) + \sum_{\beta \neq \alpha} (\phi^\beta)^2 \frac{\partial \phi^\alpha}{\partial \mathbf{n}} = \phi^\alpha N\Lambda + \left(\sum_{\beta} (\phi^\beta)^2 \right) \frac{\partial \phi^\alpha}{\partial \mathbf{n}} = \Lambda \end{aligned}$$

and thus (with $\sum_{\beta} (\phi^\beta)^2 > 0$)

$$\frac{\partial \phi^\alpha}{\partial \mathbf{n}} = \frac{\Lambda - \phi^\alpha N\Lambda}{\sum_{\beta} (\phi^\beta)^2}.$$

6.3.2 The Choice of Dynamics

In practice, the phasefield model is most commonly employed as a dynamic model for the evolution of microstructures, and not in terms of a pure minimization problem. While the choice of a phasefield functional does, within a variational framework, fix the potential equilibrium points through the solutions of the steady-state equations (6.60) resp. (6.64) above, this does not imply anything specific about the evolution of the phasefield variables, besides the very broad requirement that they should converge to a (potentially local) minimizer resp. maximizer of the phasefield functional. A simple and natural way of enforcing such a behavior is to postulate a gradient-type flow. In the Allen-Cahn case, the postulate is that of a non-conservative gradient flow based on the L^2 -gradient of \mathcal{E} , i.e. $\frac{\partial \phi}{\partial t} \sim -\frac{d\mathcal{F}_\epsilon}{d\phi}$, with individual models differing in the particular choice of proportionality to $\frac{d\mathcal{F}_\epsilon}{d\phi}$.

The model in [52] corresponds to a multiphase version of the simple scalar proportionality in Equation (6.24), i.e. to postulating

$$\tau \epsilon \frac{\partial \phi}{\partial t} = -\frac{d\mathcal{F}_\epsilon}{d\phi} - \Lambda \mathbf{e}, \quad (6.71)$$

or, in combination with the Gibbs-simplex constraint,

$$\tau \epsilon \frac{\partial \phi}{\partial t} = -\frac{d\mathcal{F}_\epsilon}{d\phi} - \Lambda \mathbf{e} + \boldsymbol{\mu}, \quad (6.72)$$

where \mathbf{e} is an N -dimensional vector of ones and $\boldsymbol{\mu}$ is subject to the same complementarity conditions as in Equation (6.64). In a more explicit form for the individual phases, this can also be written as

$$\tau \epsilon \frac{\partial \phi^\alpha}{\partial t} = \epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) - \epsilon \frac{\partial a}{\partial \phi^\alpha} - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} - \frac{\partial f}{\partial \phi^\alpha} - \Lambda, \quad \alpha = 1, \dots, N. \quad (6.73)$$

resp.

$$\tau \epsilon \frac{\partial \phi^\alpha}{\partial t} = \epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) - \epsilon \frac{\partial a}{\partial \phi^\alpha} - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} - \frac{\partial f}{\partial \phi^\alpha} - \Lambda + \mu^\alpha, \quad \alpha = 1, \dots, N. \quad (6.74)$$

together with $\sum_{\alpha=1}^N \phi^\alpha = 1$ as well as $\mu^\alpha \geq 0$ and $\mu^\alpha = 0$ if $\phi^\alpha > 0$ in the obstacle case.

One disadvantage of using a scalar-valued kinetic coefficient is that this can lead to difficulties within multiphase regions involving a mixture of highly mobile and highly immobile phases (or phase-pairings). In fact, since τ is locally the same for the evolution of all phases, an interpolation favoring the immobile interfaces may artificially slow down the evolution of the mobile ones and vice versa. This raises the question of how and to what degree one can properly interpolate between the individual mobilities such as to obtain satisfactory results.

A popular alternative, propagated in particular in [68] and used for e.g. also in [61] and [65], is again based on considering two-phase interactions, but now postulating an evolution of the form

$$\frac{\partial \phi^\alpha}{\partial t} = -\sum_{\beta \neq \alpha} m^{\alpha\beta} \left(\frac{d\mathcal{F}_\epsilon}{d\phi^\alpha} - \frac{d\mathcal{F}_\epsilon}{d\phi^\beta} \right), \quad \alpha = 1, \dots, N. \quad (6.75)$$

This corresponds to postulating that the evolution of each phase is determined as the sum of its individual interactions with all other phases, each such interaction being equipped with its

This can be recombined with Equation (6.70) - now in terms of a ϕ^α -weighted summation over the previous relation for the individual phases - to obtain

$$-N\Lambda = \sum_{\alpha} \phi^\alpha \frac{\partial \phi^\alpha}{\partial \mathbf{n}} = \frac{\sum_{\alpha} \phi^\alpha \Lambda - \sum_{\alpha} (\phi^\alpha)^2 N\Lambda}{\sum_{\beta} (\phi^\beta)^2} = \frac{\Lambda - \sum_{\alpha} (\phi^\alpha)^2 N\Lambda}{\sum_{\beta} (\phi^\beta)^2} = \frac{\Lambda}{\sum_{\beta} (\phi^\beta)^2} - N\Lambda,$$

from which it follows that Λ indeed has to be zero.

own mobility $m^{\alpha\beta} = m^{\beta\alpha}$.

This may also be written more compactly as

$$\epsilon \frac{\partial \phi}{\partial t} = -\mathbf{M} \frac{d\mathcal{F}_\epsilon}{d\phi}, \quad (\mathbf{M})_{\alpha\beta} = \begin{cases} \sum_{\gamma \neq \alpha} m^{\alpha\gamma} & \beta = \alpha, \\ -m^{\alpha\beta} & \text{else,} \end{cases} \quad (6.76)$$

and therefore amounts to replacing the scalar-valued kinetic coefficient on the left-hand side of Equation (6.71) with a symmetric **mobility-matrix** \mathbf{M} premultiplying the gradient.

It should be noted that there is no multiplier Λ for the sum-constraint in Equations (6.75) resp. (6.76). This is due to the fact that the mobility matrix \mathbf{M} above satisfies $\mathbf{M}\mathbf{e} = \mathbf{e}^T \mathbf{M} = \mathbf{0}$ and thus on the one hand $\mathbf{M} \cdot \left(\frac{d\mathcal{F}_\epsilon}{d\phi} + \Lambda \mathbf{e} \right) = \mathbf{M} \cdot \frac{d\mathcal{F}_\epsilon}{d\phi}$, i.e. applying \mathbf{M} to the gradient of the Lagrangian $\mathcal{L}(\phi) = \mathcal{F}_\epsilon(\phi) + \int_\Omega \Lambda(\mathbf{e} \cdot \phi - 1) d\mathbf{x}$ associated with the local sum constraint leads to the same result as applying \mathbf{M} directly to the gradient of \mathcal{F}_ϵ . On the other hand, as

$$\frac{\partial \sum_{\alpha=1}^N \phi^\alpha}{\partial t} = \frac{\partial \mathbf{e} \cdot \phi}{\partial t} = -\frac{1}{\epsilon} \mathbf{e} \cdot \mathbf{M} \frac{d\mathcal{F}_\epsilon}{d\phi} = \mathbf{0},$$

there is no need to explicitly introduce the multiplier Λ since the phasefield will remain consistent with the sum-constraint provided the initial values were so.

Remark 71. This is analogous in nature to the choice of the mobility matrix \mathbf{L} in [52] for the concentration- and energy-evolution, which, in addition to the multiplier Λ in the phasefield equation would a priori also require another one Λ_c for the concentration equation due to the constraint $\sum_{i=1}^K c_i = 1$. As \mathbf{L} is chosen to be only a positive semi-definite matrix satisfying $\sum_{i=1}^K L_{ij} = \sum_{j=1}^K L_{ij} = 0$ for all rows resp. columns, the multiplier for the concentration drops out of the resulting equations while the sum-constraint nevertheless remains satisfied provided it was so initially. \diamond

This is in contrast to the multipliers $\boldsymbol{\mu}$ arising due to the inequality constraints in the obstacle case, which still need to be considered explicitly such that the evolution equation has to be modified to (for the representation in (6.2) of the Gibbs-simplex)

$$\epsilon \frac{\partial \phi}{\partial t} = -\mathbf{M} \left(\frac{d\mathcal{F}_\epsilon}{d\phi} - \boldsymbol{\mu} \right), \quad \mu^\alpha \phi^\alpha = 0, \quad \mu^\alpha \geq 0, \quad \alpha = 1, \dots, N. \quad (6.77)$$

Even though this formulation still has the apparent advantage of containing one multiplier less than the obstacle-version of Equation (6.71) (the sum-constraint remaining satisfied regardless of the choice of $\boldsymbol{\mu}$ as well), this simply shifts the indirect interaction of each non-zero multiplier μ^α with the other phases through the sum-constraint to a more explicit one in terms of the full matrix \mathbf{M} .

Remark 72. Note that one could in principle also rewrite Equation (6.71) resp. its obstacle-analogon in a similar manner as⁴⁸

$$\frac{\partial \phi}{\partial t} = -\frac{1}{\tau\epsilon} \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \otimes \mathbf{e} \right) \frac{d\mathcal{F}_\epsilon}{d\phi} \quad \text{resp.} \quad \frac{\partial \phi}{\partial t} = -\frac{1}{\tau\epsilon} \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \otimes \mathbf{e} \right) \left(\frac{d\mathcal{F}_\epsilon}{d\phi} - \boldsymbol{\mu} \right),$$

where the Lagrange-multiplier Λ has been eliminated through the use of the ‘‘mobility matrix’’ $\mathbf{M} = \frac{1}{\tau\epsilon} \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \otimes \mathbf{e} \right)$ corresponding to a scalar multiple of the Euclidian projection operator onto the subspace of zero-average vectors. As the diagonal entries of this projection matrix are given by $\frac{1}{\tau\epsilon} \left(1 - \frac{1}{N} \right) = \frac{N-1}{N\tau\epsilon} = -\sum_{\beta \neq \alpha} \left(-\frac{1}{N\tau\epsilon} \right)$, this is just a particular simple form of the one above obtained by setting $m^{\alpha\beta} = -\frac{1}{N\tau\epsilon} \quad \forall \alpha \neq \beta$.

⁴⁸This corresponds to a (partially in the obstacle case) ‘‘projected form’’ of the equation.

It is therefore not surprising that both approaches actually share many of the same features and difficulties, the most notable difference being that the particularly simple form in the latter case often allows for more “explicit calculations”.

One point where the more complex form of the mobility matrix \mathbf{M} leads to a notable complication is when relying on a projection-based descent algorithm. In fact, the modified interaction of the gradients \mathbf{g} through \mathbf{M} also has to be taken into account during this projection operation (in line with the role of $\boldsymbol{\mu}$ in Equation (6.77)), and entails the necessity of using a more complex \mathbf{M} -weighted projection operator instead of the simpler Euclidian one. This projection, in particular some of its algorithmic implications will be discussed in Section 6.3.3.

It can also be noted that in the two-phase case, the matrix $\frac{1}{\tau\epsilon}(\mathbf{I}-\mathbf{e}\otimes\mathbf{e})$ reduces to $\frac{1}{2\tau\epsilon}\begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$, and both approaches are therefore clearly equivalent if m^{12} is chosen as $\frac{1}{2\tau}$. \diamond

Remark 73. The two approaches outlined above are clearly not the only ones possible. In particular, one can in principle use any spd mobility matrix \mathbf{M} provided one is willing to maintain the explicit multiplier Λ . Similarly, one could also choose to replace the use of a mobility matrix by that of a matrix $\boldsymbol{\tau}$ of kinetic coefficients and postulating $\boldsymbol{\tau}\frac{\partial\phi}{\partial t} = \frac{d\mathcal{E}}{d\phi} - \Lambda$ (this of course being equivalent to the previous point in the spd case by setting $\mathbf{M} = \boldsymbol{\tau}^{-1}$ or vice versa)⁴⁹. The latter may be potentially useful for obtaining a somewhat “intermediate” model between the scalar kinetic coefficient and the use of a full mobility matrix if applied in the form

$$\tau^\alpha \epsilon \frac{\partial\phi^\alpha}{\partial t} = -\frac{d\mathcal{F}_\epsilon}{d\phi^\alpha} - \Lambda, \quad \alpha = 1, \dots, N$$

i.e. by using a simple diagonal matrix $\boldsymbol{\tau} = \text{diag}((\tau^\alpha)_{1 \leq \alpha \leq N})$ allowing for some additional freedom in specifying the dynamics while remaining in an “almost explicit” form due to the diagonal structure of $\boldsymbol{\tau}$. \diamond

6.3.3 Some Numerical and Algorithmic Considerations

A first important point to observe from a structural point of view is that the derivative contribution from $w(\phi)$ now always has a nonlinear contribution due to the triple-phase terms, even in the previously linear obstacle case. Depending on the size of the penalty parameters $\gamma^{\alpha\beta\delta}$, this contribution may be relatively large and can thus potentially also lead to additional numerical difficulties. Similarly, due to the formulation of the surface energy densities in terms of the (nonlinear) $\mathbf{q}^{\alpha\beta}$, the derivative contribution due to a will also always be nonlinear both in ϕ and $\nabla\phi$, even in the previously linear isotropic setting.

An additional complication, in particular in the obstacle case, arises due to the nature of the admissible set. Whereas in the two-phase case, the sum-constraint is easily absorbed into a reduced formulation and thus leads to no actual constraint in the double-well case and a simple box-constraint in the double-obstacle case, it is more common when dealing with multiphase problems to remain within a non-reduced formulation. This requires maintaining the sum-constraint as an explicit constraint. While this is generally easily dealt with in the well case, the combination with the positivity constraint is a somewhat trickier issue. In particular, there is no simple explicit formula even for the simplest (Euclidian) projector onto the Gibbs-simplex, as the combinatoric nature of the projection operation is more pronounced due to the higher complexity of the constraint set (regardless of whether one uses a reduced formulation or not).

⁴⁹In fact, if one maintains the explicit presence of the Lagrange multiplier Λ , one could also use a spsd $\boldsymbol{\tau}$ -matrix whose kernel (and, by symmetry, cokernel) coincides with $\text{Span}\{\mathbf{e}\}$, even though this seems of somewhat limited practical interest. While $\boldsymbol{\tau}$ is then not invertible, the same choice of Λ (up to a prefactor of $\tau\epsilon$) as for Equation (6.71) ensures consistency (i.e. orthogonality to the constants) of the left- and right-hand side, while the indeterminacy in $\frac{\partial\phi}{\partial t}$ is eliminated by the presence of the sum-constraint.

Remark 74. Even though it is in principle possible to use a reduced formulation in the multiphase setting by eliminating one of the N phasefields in terms of the remaining $N - 1$, any potential benefits of doing so would primarily be restricted to the multi-well potential. The central difficulty in the multi-obstacle case is not the sum-constraint, but its coupling with the bounds on the ϕ -values. As this implies the risk that the phase which was eliminated in favor of the others may actually not be allowed to change according to the changes of the other ones, this does not in general avoid the issue of the sum-constraint. This is in contrast to the two-phase case, where any one phase trying to move outside one of the 0-1-bound is equivalent to the other one trying to move outside the opposite one. \diamond

Before continuing the discussion on the challenges associated with “true” multi-phasefield problems with $N > 2$, it is worthwhile to first take a closer look at how the formulation above compares with the simpler description in terms of a single phasefield-variable ϕ from the previous section. While it is clear that they are in principle equivalent, the redundant representation in terms of $\phi^1 := \phi$ and $\phi^2 = 1 - \phi^1 = 1 - \phi$ leads to a number of “technicalities” which should be distinguished from the actual difficulties associated with the situation for $N > 2$.

Comparison of the “Single-Phase” to the “Two-Phase” Formulation

If one **explicitly** enforces the sum-constraint by reexpressing e.g. ϕ^2 as a function of ϕ^1 , i.e. with $\phi^2 = 1 - \phi^1$, $\nabla\phi^2 = -\nabla\phi^1$ and thus also $\mathbf{q}^{12} = -\nabla\phi^1$, it is easily seen that the values of the a - and w -term (not the functions though!) reduce to the same expressions that would be obtained from Equation (6.16). If the reduced functions in terms of $\phi^2 = \phi^2(\phi^1)$ are denoted with a hat-symbol, i.e. for example $\hat{w}(\phi^1) = w(\phi^1, \phi^2(\phi^1))$, as a simple consequence of the chain rule and of

$$\frac{\partial \mathbf{q}^{12}}{\partial \phi^1} = \nabla\phi^2, \quad \frac{\partial \mathbf{q}^{12}}{\partial \phi^2} = -\nabla\phi^1, \quad \frac{\partial \mathbf{q}^{12}}{\partial \nabla\phi^1} = -\phi^2 \mathbf{I}, \quad \frac{\partial \mathbf{q}^{12}}{\partial \nabla\phi^2} = \phi^1 \mathbf{I}$$

the respective derivatives in the phasefield equation are given by

$$\begin{aligned} \frac{1}{\gamma^{12}} \frac{\delta \hat{a}(\phi^1)}{\delta \phi^1} &= \frac{1}{\gamma^{12}} \frac{\delta a}{\delta \phi^1} - \frac{\delta a}{\delta \phi^2} = -\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^1} - \frac{\partial a}{\partial \nabla \phi^2} \right) + \left(\frac{\partial a}{\partial \phi^1} - \frac{\partial a}{\partial \phi^2} \right) \\ &= -\nabla \cdot \left(- \underbrace{\phi^2}_{=1-\phi^1} \mathbf{q}^{12} + \phi^1 \mathbf{q}^{12} \right) + \mathbf{q}^{12} \cdot \left(\underbrace{\nabla \phi^2}_{=-\nabla \phi^1} - (-\nabla \phi^1) \right) = \nabla \cdot \mathbf{q}^{12} = -\Delta \phi^1 \end{aligned}$$

and

$$\begin{aligned} \frac{\partial \hat{w}_{dw}}{\partial \phi^1} &= \frac{\partial w_{dw}}{\partial \phi^1} - \frac{\partial w_{dw}}{\partial \phi^2} = 18\gamma^{12}(\phi^1(\phi^2)^2 - (\phi^1)^2\phi^2) = 18\gamma^{12}\phi^1\phi^2(\phi^2 - \phi^1) \\ &= 18\gamma^{12}\phi^1(1 - \phi^1)(1 - 2\phi^1) \end{aligned}$$

for the double-well potential and

$$\frac{\partial \hat{w}_{do}}{\partial \phi^1} = \frac{\partial w_{do}}{\partial \phi^1} - \frac{\partial w_{do}}{\partial \phi^2} = \frac{16}{\pi^2}\gamma^{12}(\phi^2 - \phi^1) = \frac{16}{\pi^2}\gamma^{12}(1 - 2\phi^1)$$

for the double-obstacle potential respectively. In addition, if \tilde{h} is based on any of the “standard” two-phase interpolation functions and $\phi^1 + \phi^2 = 1$, and thus also $\tilde{h}^1(\phi) + \tilde{h}^2(\phi) = 1$, then the denominator drops out and the values of the derivatives reduce to $\frac{\partial \hat{h}^1}{\partial \phi^1} = \tilde{h}(\phi^2)\tilde{h}'(\phi^1)$ resp. $\frac{\partial \hat{h}^1}{\partial \phi^2} = -\tilde{h}(\phi^1)\tilde{h}'(\phi^2)$. As one further has $\tilde{h}'(\phi^2) = -\tilde{h}'(\phi^1)$, the derivative of the f -contribution

becomes

$$\begin{aligned}
\frac{\partial \hat{f}(\phi^1)}{\partial \phi^1} &= \frac{\partial f}{\partial \phi^1}((\phi^1, \phi^2(\phi^1))) = \frac{\partial f}{\partial \phi^1} - \frac{\partial f}{\partial \phi^2} = f^1 \left(\frac{\partial h^1}{\partial \phi^1} - \frac{\partial h^1}{\partial \phi^2} \right) + f^2 \left(\frac{\partial h^2}{\partial \phi^1} - \frac{\partial h^2}{\partial \phi^2} \right) \\
&= f^1 (\tilde{h}(\phi^2) \tilde{h}'(\phi^1) - \tilde{h}(\phi^1) \tilde{h}'(\phi^2)) + f^2 (\tilde{h}(\phi^1) \tilde{h}'(\phi^2) - \tilde{h}(\phi^2) \tilde{h}'(\phi^1)) \\
&= (f^1 - f^2) \left(\underbrace{\tilde{h}(\phi^2)}_{=1-\tilde{h}(\phi^1)} \tilde{h}'(\phi^1) - \tilde{h}(\phi^1) \underbrace{\tilde{h}'(\phi^2)}_{=-\tilde{h}'(\phi^1)} \right) = (f^1 - f^2) \tilde{h}'(\phi^1),
\end{aligned}$$

showing that this term as well reduces to the previous expression.

This is of course not particularly surprising as the multi-phasefield formulation is naturally intended as a generalization of the two-phase case. The simple calculations above nevertheless highlight two important points:

1. Whereas the original (isotropic) “one-phase” formulation is linear in the second-order term and only involves at most low-order polynomial expressions in the derivatives of the w - and/or f -term, its “two-phase” counterpart is a priori also nonlinear in the second order term and involves much more unpleasant nonlinearities (fractions over squares of the original h -function) in the f -term.
2. If the problem is not from the outset reduced to the previous version by using either $\phi^2 = \phi^2(\phi^1)$ or $\phi^1 = \phi^1(\phi^2)$, the cancellations occurring in the reduced formulation will instead have to be enforced through an external Lagrange multiplier λ for the sum-constraint $\phi^1 + \phi^2 = 1$. This, by itself, poses no particular problem. Nevertheless, if one additionally chooses to (or has to) enforce the constraints $0 \leq \phi^1, \phi^2 \leq 1$, a direct extension of the procedure applied in the two-phase case using separate multipliers $\mu^{1/2, \pm}$ for the lower and upper constraints entails some additional technical difficulties. Due to the sum constraint, there will be either no active inequality constraint or **both** equality constraints will be active at the same time. In the former case, the Lagrange multiplier λ for the sum-constraint will, as expected, just equal the average of both derivatives. In the latter case though, there are a total of three multipliers (for example λ , $\mu^{1,-}$ and $\mu^{2,+}$ if one has $\phi^1 = 0$, $\phi^2 = 1$) to be determined based on only two actual unknowns ϕ^1 and ϕ^2 and one is faced with an underdetermined problem. Barring the exceptional case in which both constraints are only weakly active (i.e. if both phases would remain in 0 resp. 1 even without the box-constraints), this necessarily results in a non-unique (though bounded) set of Lagrange multipliers.

This is not a problem per se, as there is no real need to have uniqueness for the Lagrange multipliers. In particular, this indeterminacy can be dealt with through an appropriate selection mechanism and the resulting phasefield-values will be the same regardless of the particular choice made. Depending on the way λ and the $\mu^{\alpha, \pm}$ are determined, this nevertheless may require some care when e.g. working with direct solvers for the associated matrices.

Remark 75. A similar issue with respect to the multipliers arises also if one e.g. imposes an additional constraint $\int_{\Omega} \phi^{\alpha} = V^{\alpha}$ on the phase-volumes. These a priori two constraints require an additional set of Lagrange multipliers (χ^1, χ^2) , leading to the modification

$$\tau \epsilon \frac{\partial \phi^{\alpha}}{\partial t} = -g^{\alpha} - \lambda - \chi^{\alpha}, \quad \alpha = 1, 2,$$

of the right-hand sides (or, when using the slightly modified approach in [53] and [29], of the form $\sum_{\beta=1}^2 \chi^{\beta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}$) and the appropriate modification

$$\lambda(\mathbf{x}, t) = \frac{1}{2} \left(-g^1(\mathbf{x}, t) - g^2(\mathbf{x}, t) - (\chi^1 + \chi^2) \right) \quad (6.78)$$

of λ . There is again a redundancy in this case as it would be sufficient to fix just one of the volumes, e.g. that of the first phase, and then let λ “handle the rest”. In particular with respect to splitting-type approaches for the above system, it is highly preferable to maintain, at least formally, both χ^1 and χ^2 as separate multipliers as this avoids enforcing the global constraint on the other phase indirectly through a multiplier which is intuitively associated with the local constraint. A natural selection criterion for eliminating the indeterminacy in this case is given by choosing as in [53] χ^1 and χ^2 such $\chi^1 + \chi^2 = 0$, since, by Equation (6.78), this leaves λ unaffected by the additional constraint. \diamond

Some Choices for the Spatial Discretization

The presence of more than two phases has essentially no impact on the simplest discretization of the local terms due to the bulk potential and the driving forces within the cell-centered scheme used in the **Pace3D**-framework. In contrast to the pure two-phase setting, the derivative of the gradient energy density in the simplest isotropic case cannot be reduced to a Laplacian anymore as soon as there are more than two phases present, and one does thus not have an “obvious” choice of discretization through the standard second-order difference quotient anymore⁵⁰.

Based on the divergence-type nature of the expression $\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right)$, a very convenient discretization for this term is to express it as the differences of fluxes over the cell “faces”, i.e. in the simplest one-dimensional setting as

$$\epsilon \frac{d}{dx} \left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right) \approx \frac{1}{\Delta x} \left(\left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right)_{i+\frac{1}{2}} - \left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right)_{i-\frac{1}{2}} \right),$$

where the $\left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right)_{j+\frac{1}{2}}$ correspond to appropriate approximations of this derivative on these

faces and the prime indicates a derivative with respect to x . As $\left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right)$ is in the isotropic case (this being the only reasonable choice in 1D) given by

$$\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} = -2 \left(\sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \phi^\beta q^{\alpha\beta} \right)$$

and $q^{\alpha\beta} \left(\phi, \frac{d\phi}{dx} \right) = \phi^\alpha \frac{d\phi^\beta}{dx} - \phi^\beta \frac{d\phi^\alpha}{dx}$, the most obvious choice of discretization for this expression on the faces is obtained by using a short one-cell stencil for the first-order derivatives and the average of the neighboring ϕ -values for the values of the phasefield itself. Defining $\phi_{i+\frac{1}{2}} = \frac{1}{2}(\phi_i + \phi_{i+1})$, this leads to the natural approximations

$$q_{i+\frac{1}{2}}^{\alpha\beta} = \phi_{i+\frac{1}{2}}^\alpha \frac{\phi_{i+1}^\beta - \phi_i^\beta}{\Delta x} - \phi_{i+\frac{1}{2}}^\beta \frac{\phi_{i+1}^\alpha - \phi_i^\alpha}{\Delta x} \quad (6.79)$$

and

$$\epsilon \left(\frac{d}{dx} \left(\frac{\partial a}{\partial \left(\frac{d\phi^\alpha}{dx} \right)} \right) \right)_i \approx -2\epsilon \frac{1}{\Delta x} \left[\left(\sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \phi_{i+\frac{1}{2}}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} \right) - \left(\sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \phi_{i-\frac{1}{2}}^\beta q_{i-\frac{1}{2}}^{\alpha\beta} \right) \right]. \quad (6.80)$$

In contrast, there are actually two quite natural discretizations for the second contribution from the term $\epsilon \frac{\partial a}{\partial \phi^\alpha}$, which is, in the isotropic case, given by $\epsilon \frac{\partial a}{\partial \phi^\alpha} = 2\epsilon \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} q^{\alpha\beta} \frac{d\phi^\beta}{dx}$. The first and most intuitive one is to evaluate it in a cell-centered fashion, i.e. by using the approximation

$$\left(\epsilon \frac{\partial a}{\partial \phi^\alpha} \right)_i \approx 2\epsilon \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \left(\phi_i^\alpha \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} - \phi_i^\beta \frac{\phi_{i+1}^\alpha - \phi_{i-1}^\alpha}{2\Delta x} \right) \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} \quad (6.81)$$

⁵⁰Using higher-order stencils is of course in principle possible, but, in the obstacle case, unlikely to lead to any major improvement by the discussion in Section 6.2.3.

through central gradients only. An alternative is to try to reuse the already calculated gradients of ϕ and the $q^{\alpha\beta}$ on the cell faces and thus to express $\epsilon \frac{\partial a}{\partial \phi^\alpha}$ as the average of its approximation on the left and right face as

$$\left(\epsilon \frac{\partial a}{\partial \phi^\alpha} \right)_i \approx 2\epsilon \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \frac{1}{2} \left(q_{i-\frac{1}{2}}^{\alpha\beta} \frac{\phi_i^\beta - \phi_{i-1}^\beta}{\Delta x} + q_{i+\frac{1}{2}}^{\alpha\beta} \frac{\phi_{i+1}^\beta - \phi_i^\beta}{\Delta x} \right). \quad (6.82)$$

The use of this - a priori slightly less intuitive - expression as compared to the one in Equation (6.81) is two-fold. Firstly, it can indeed partially rely upon the same face-centered quantities $q_{j+\frac{1}{2}}^{\alpha\beta}$ already required for the evaluation of the divergence-term in Equation (6.80), but at the price of introducing the same spatial buffering requirement as for the fluxes over the cell faces if one wants to avoid recalculating the same values twice. Secondly, the resulting discrete phasefield equation based on the combined discretization of the gradient energy contributions as in equations (6.80) and (6.82) together with the appropriate isolating boundary conditions corresponds itself to the derivative of a discrete energy “functional” $\mathcal{F}_{\epsilon, \Delta x}$ given by

$$\mathcal{F}_{\epsilon, \Delta x}(\phi) = \sum_i \frac{1}{2} \epsilon \left(a(\phi_{i-\frac{1}{2}}, \frac{\phi_i - \phi_{i-1}}{\Delta x}) + a(\phi_{i+\frac{1}{2}}, \frac{\phi_{i+1} - \phi_i}{\Delta x}) \right) + \frac{1}{\epsilon} w(\phi_i) + f(\phi_i), \quad (6.83)$$

where $a(\phi_{j+\frac{1}{2}}, \frac{\phi_{j+1} - \phi_j}{\Delta x})$ abbreviates the summation $\sum_\alpha \sum_{\beta > \alpha} |q_{j+\frac{1}{2}}^{\alpha\beta}|^2$ with $q_{j+\frac{1}{2}}^{\alpha\beta}$ defined as in Equation (6.79).

This is obvious for the local terms arising due to w and f . For seeing that this is also true with respect to the spatial terms in the gradient energy contribution, it is convenient to first rewrite some of the expressions above.

Firstly, it is easy to verify that the expression for $q^{\alpha\beta}$ in Equation (6.79) can be simplified to

$$q_{i+\frac{1}{2}}^{\alpha\beta} = \frac{\phi_i^\alpha \phi_{i+1}^\beta - \phi_i^\beta \phi_{i+1}^\alpha}{\Delta x}. \quad (6.84)$$

Secondly, the combined contributions from the divergence term and the derivative of a with respect to ϕ^α can be summarized to

$$\begin{aligned} & -2\epsilon \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \left\{ \frac{1}{\Delta x} \left[\phi_{i+\frac{1}{2}}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} - \phi_{i-\frac{1}{2}}^\beta q_{i-\frac{1}{2}}^{\alpha\beta} \right] + \frac{1}{2} \left(q_{i-\frac{1}{2}}^{\alpha\beta} \frac{\phi_i^\beta - \phi_{i-1}^\beta}{\Delta x} + q_{i+\frac{1}{2}}^{\alpha\beta} \frac{\phi_{i+1}^\beta - \phi_i^\beta}{\Delta x} \right) \right\} \\ &= -2 \frac{\epsilon}{\Delta x} \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \left[q_{i+\frac{1}{2}}^{\alpha\beta} \left(\phi_{i+\frac{1}{2}}^\beta + \frac{1}{2} (\phi_{i+1}^\beta - \phi_i^\beta) \right) - q_{i-\frac{1}{2}}^{\alpha\beta} \left(\phi_{i-\frac{1}{2}}^\beta - \frac{1}{2} (\phi_i^\beta - \phi_{i-1}^\beta) \right) \right] \\ &= -2 \frac{\epsilon}{\Delta x} \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \left[\phi_{i+1}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} - \phi_{i-1}^\beta q_{i-\frac{1}{2}}^{\alpha\beta} \right]. \end{aligned}$$

Differentiating the total gradient energy contribution in Equation (6.83) with respect to a particular ϕ -value ϕ_i^α , it is easy to see that it suffices to verify that for each phase-pairing individually, one has

$$-\frac{d}{d\phi_i^\alpha} \sum_j \frac{1}{2} \left(|q_{j-\frac{1}{2}}^{\alpha\beta}|^2 + |q_{j+\frac{1}{2}}^{\alpha\beta}|^2 \right) = -2 \frac{\epsilon}{\Delta x} \phi_{i+1}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} - \phi_{i-1}^\beta q_{i-\frac{1}{2}}^{\alpha\beta}. \quad (6.85)$$

Since ϕ_i^α is only involved in two terms in this summation, namely the ones for $j = i$, it follows that this is the same as

$$-\frac{d}{d\phi_i^\alpha} \frac{1}{2} \left(|q_{i-\frac{1}{2}}^{\alpha\beta}|^2 + |q_{i+\frac{1}{2}}^{\alpha\beta}|^2 \right) = -2 \left(q_{i-\frac{1}{2}}^{\alpha\beta} \frac{dq_{i-\frac{1}{2}}^{\alpha\beta}}{d\phi_i^\alpha} + q_{i+\frac{1}{2}}^{\alpha\beta} \frac{dq_{i+\frac{1}{2}}^{\alpha\beta}}{d\phi_i^\alpha} \right),$$

and thus together with $\frac{dq_{i+\frac{1}{2}}^{\alpha\beta}}{d\phi_i^\alpha} = \frac{\phi_{i+1}^\beta}{\Delta x}$ and $\frac{dq_{i-\frac{1}{2}}^{\alpha\beta}}{d\phi_i^\alpha} = \frac{d}{d\phi_i^\alpha} \frac{\phi_{i-1}^\alpha \phi_i^\beta - \phi_{i-1}^\beta \phi_i^\alpha}{\Delta x} = -\frac{\phi_{i-1}^\beta}{\Delta x}$ based on the simplified expression in Equation (6.84) that Equation (6.85) is indeed satisfied.

In contrast, a similar property does not hold when using a discretization of the gradient energy distributions using equations (6.80) and (6.82). In fact, the existence of an underlying potential would, by Schwarz's theorem, require the symmetry of the second derivatives with respect to all unknowns. Even though it is clear that both contributions are symmetric with respect to the phase-indices $\alpha \in \{1, 2, \dots, N\}$, this is not the case with respect to the ϕ -values at different spatial positions. Combining equations (6.80) and (6.82), the total gradient energy contribution in a cell i is given by

$$-\tilde{g}_i^\alpha = -2\epsilon \sum_{\beta \neq \alpha} \gamma^{\alpha\beta} \left\{ \frac{1}{\Delta x} \left[\phi_{i+\frac{1}{2}}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} - \phi_{i-\frac{1}{2}}^\beta q_{i-\frac{1}{2}}^{\alpha\beta} \right] + \left(\phi_i^\alpha \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} - \phi_i^\beta \frac{\phi_{i+1}^\alpha - \phi_{i-1}^\alpha}{2\Delta x} \right) \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} \right\}.$$

In order to be compatible with the symmetry of the second derivatives, one would therefore have to satisfy $\frac{\partial \tilde{g}_i^\alpha}{\partial \phi_j^\alpha} = \frac{\partial \tilde{g}_j^\alpha}{\partial \phi_i^\alpha}$ for any combination of i and j , or, restricting the attention to two neighboring cells $\frac{\partial \tilde{g}_i^\alpha}{\partial \phi_{i+1}^\alpha} = \frac{\partial \tilde{g}_{i+1}^\alpha}{\partial \phi_i^\alpha}$. By the symmetry with respect to the phases, it is furthermore again clear that it suffices to check this for the contribution of a single phase-pairing, and is thus a matter of comparing

$$\begin{aligned} & \frac{\partial}{\partial \phi_{i+1}^\alpha} \left\{ \frac{1}{\Delta x} \left[\phi_{i+\frac{1}{2}}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} - \phi_{i-\frac{1}{2}}^\beta q_{i-\frac{1}{2}}^{\alpha\beta} \right] + \left(\phi_i^\alpha \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} - \phi_i^\beta \frac{\phi_{i+1}^\alpha - \phi_{i-1}^\alpha}{2\Delta x} \right) \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} \right\} \\ &= \frac{1}{\Delta x} \phi_{i+\frac{1}{2}}^\beta \frac{\partial q_{i+\frac{1}{2}}^{\alpha\beta}}{\partial \phi_{i+1}^\alpha} - \phi_i^\beta \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{2\Delta x} \frac{\partial}{\partial \phi_{i+1}^\alpha} \frac{\phi_{i+1}^\alpha - \phi_{i-1}^\alpha}{2\Delta x} = \frac{1}{\Delta x} \phi_{i+\frac{1}{2}}^\beta \frac{\partial q_{i+\frac{1}{2}}^{\alpha\beta}}{\partial \phi_{i+1}^\alpha} - \phi_i^\beta \frac{\phi_{i+1}^\beta - \phi_{i-1}^\beta}{4(\Delta x)^2} \end{aligned}$$

and

$$\begin{aligned} & \frac{\partial}{\partial \phi_i^\alpha} \left\{ \frac{1}{\Delta x} \left[\phi_{i+\frac{3}{2}}^\beta q_{i+\frac{3}{2}}^{\alpha\beta} - \phi_{i+\frac{1}{2}}^\beta q_{i+\frac{1}{2}}^{\alpha\beta} \right] + \left(\phi_{i+1}^\alpha \frac{\phi_{i+2}^\beta - \phi_i^\beta}{2\Delta x} - \phi_{i+1}^\beta \frac{\phi_{i+2}^\alpha - \phi_i^\alpha}{2\Delta x} \right) \frac{\phi_{i+2}^\beta - \phi_i^\beta}{2\Delta x} \right\} \\ &= -\frac{1}{\Delta x} \phi_{i+\frac{1}{2}}^\beta \frac{\partial q_{i+\frac{1}{2}}^{\alpha\beta}}{\partial \phi_i^\alpha} - \phi_{i+1}^\beta \frac{\phi_{i+2}^\beta - \phi_i^\beta}{2\Delta x} \frac{\partial}{\partial \phi_i^\alpha} \frac{\phi_{i+2}^\alpha - \phi_i^\alpha}{2\Delta x} = -\frac{1}{\Delta x} \phi_{i+\frac{1}{2}}^\beta \frac{\partial q_{i+\frac{1}{2}}^{\alpha\beta}}{\partial \phi_i^\alpha} + \phi_{i+1}^\beta \frac{\phi_{i+2}^\beta - \phi_i^\beta}{4(\Delta x)^2}. \end{aligned}$$

Since $q_{i+\frac{1}{2}}^{\alpha\beta}$ depends neither on ϕ_{i-2}^β nor on ϕ_{i+2}^β , it is obvious that this equality cannot hold in general, and that there is thus no potential underlying this discretization.

Remark 76. This lack of variational interpretation in the discrete case is not necessarily an issue and does not automatically favor the other discretization. Instead, this choice (as well as that of any other discretization), also depends heavily on the interplay of various factors.

On the one hand, both discretizations are perfectly legitimate and differ only by a second-order discretization error. The results obtained by the two discretizations are therefore not expected to change beyond the level of the other errors inherent in the discretization anyway⁵¹.

On the other hand, it does not have any real effect on the characterization of the ‘‘multipliers’’ in Equation (6.64) (resp. its simpler well-version) or their alternative introduction through a projection-based algorithm below. In fact, even though these are then not strictly speaking Lagrange multipliers in the sense of an underlying Lagrangian⁵², they are still induced in the same form by the constraints due to the more general geometrical characterization through orthogonality and normal cone conditions for the admissible variations discussed in Sections 4 and 4.3.

It therefore often makes sense to prefer one or the other based on more practical considerations.

⁵¹This can also be verified numerically, with both solutions differing only very slightly.

⁵²Recall that the two primary motivations for the Lagrangians are the formal simplicity with which they allow to derive necessary conditions for constrained variational problems as well as any potential benefits obtained through a saddle-point structure. Neither of these is used when arguing based on Equation (6.59).

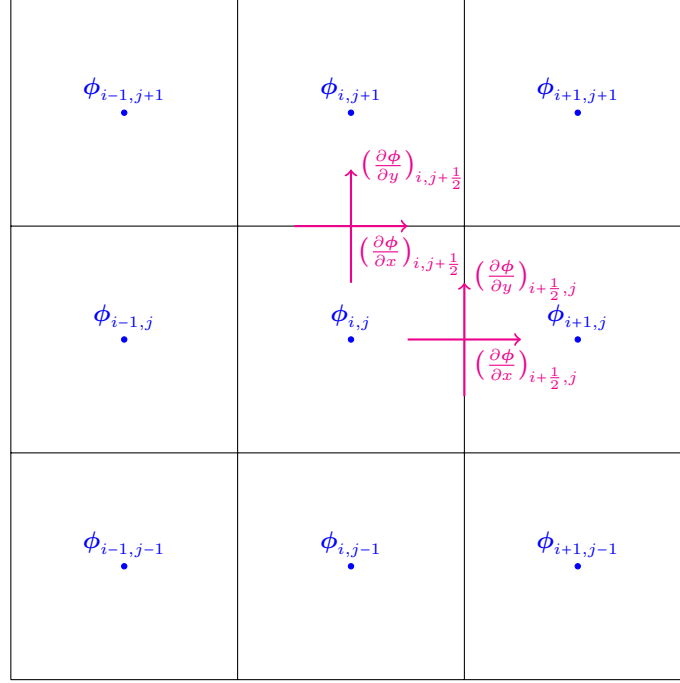


Figure 6.6: Layout of the ϕ -values and required gradients on the upper cell faces for the evaluation of the gradient energy contribution in the cell (i, j) .

In particular, when using an explicit time-discretization which is typically quite “forgiving” to moderate perturbations in the equations, a much more relevant criterion is often the difference in run-time obtained for both discretizations due to the usually very high number of time-steps inherited by the stability restrictions. In contrast, other more complex schemes can often be relatively sensible to even small perturbations in the equations. For example, the common usage of line-search or trust-region algorithms in constrained optimization problems depends fundamentally on the existence of a potential, and the use of an alternative merit-function which is not actually being minimized can thus be an issue. Finally, there is also the very “mundane” issue of the effort and amount of code required for the implementation of the equations. \diamond

In the isotropic case considered above, both discretizations are easily extended to the higher-dimensional setting by “splitting” the contributions of $\mathbf{q}^{\alpha\beta}$ onto the faces in the respective spatial direction, i.e. in two dimensions by associating the contributions from the x -component of $\mathbf{q}^{\alpha\beta}$ on the left and right cell faces, and those due to the y -component on the bottom and top faces. With some further modifications, they can also be extended to an anisotropic setting. In this respect, the primary difficulty faced by a cell-centered discretization in combination with fluxes discretized on the faces is that these in general then depend on the whole vector $\mathbf{q}^{\alpha\beta}$ resp. its orientation $\frac{\mathbf{q}^{\alpha\beta}}{|\mathbf{q}^{\alpha\beta}|}$. Whereas the layout is very convenient for calculating e.g. the x -component of $\nabla\phi$ at the face-centers in the x -direction, it is much less so for calculating the y -component at these same points as illustrated in Figure 6.6. While $(\frac{\partial\phi}{\partial x})_{i+\frac{1}{2},j}$ and $(\frac{\partial\phi}{\partial y})_{i,j+\frac{1}{2}}$ can as before be evaluated using a short difference as $\frac{\phi_{i+1,j}-\phi_{i,j}}{\Delta x}$ and $\frac{\phi_{i,j+1}-\phi_{i,j}}{\Delta y}$, the evaluation of the respective other component is most naturally done using a broader “averaged” stencil as $(\frac{\partial\phi}{\partial y})_{i+\frac{1}{2},j} \approx \frac{(\phi_{i,j+1}+\phi_{i+1,j+1})-(\phi_{i,j-1}+\phi_{i+1,j-1})}{4\Delta y}$ and $(\frac{\partial\phi}{\partial x})_{i,j+\frac{1}{2}} \approx \frac{(\phi_{i+1,j}+\phi_{i+1,j+1})-(\phi_{i-1,j}+\phi_{i-1,j+1})}{4\Delta x}$.

Based on this stencil and the same simple averaging of the ϕ -values as before for approxim-

ing $\phi_{i+\frac{1}{2},j}$ and $\phi_{i,j+\frac{1}{2}}$, one can then recover the full $\mathbf{q}^{\alpha\beta}$ -vector on each of the faces. Based on Equation (6.67), the two-dimensional and potentially anisotropic extensions of the divergence-contribution in Equation (6.80) is thus given by

$$\begin{aligned} \epsilon \left(\nabla \cdot \left(\frac{\partial a}{\partial(\nabla\phi^\alpha)} \right) \right)_i &\approx -2\epsilon \sum_{\beta \neq \alpha} \left\{ \frac{1}{\Delta x} \left[\phi_{i+\frac{1}{2},j}^\beta \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i+\frac{1}{2},j}^{\alpha\beta}) - \phi_{i-\frac{1}{2},j}^\beta \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i-\frac{1}{2},j}^{\alpha\beta}) \right] \right. \\ &\quad \left. + \frac{1}{\Delta y} \left[\phi_{i,j+\frac{1}{2}}^\beta \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i,j+\frac{1}{2}}^{\alpha\beta}) - \phi_{i,j-\frac{1}{2}}^\beta \frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i,j-\frac{1}{2}}^{\alpha\beta}) \right] \right\}. \end{aligned} \quad (6.86)$$

Similarly, one can extend the additional contribution by $\frac{\partial a}{\partial\phi^\alpha}$ in (6.66) in a cell-centered fashion as in Equation (6.81), i.e. by calculating $\nabla\phi_{i,j}$ through $\left(\frac{\partial\phi}{\partial x}\right)_{i,j} \approx \frac{\phi_{i+1,j}-\phi_{i-1,j}}{2\Delta x}$ and $\left(\frac{\partial\phi}{\partial y}\right)_{i,j} \approx \frac{\phi_{i,j+1}-\phi_{i,j-1}}{2\Delta x}$ and combining these with $\phi_{i,j}$ to obtain $\mathbf{q}_{i,j}^{\alpha\beta}$, thus leading to the total contribution

$$\left(\epsilon \frac{\partial a}{\partial\phi^\alpha} \right)_{i,j} \approx 2\epsilon \sum_{\beta \neq \alpha} \left(\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i,j}^{\alpha\beta}) \right) \cdot \begin{pmatrix} \frac{\phi_{i+1,j}^\beta - \phi_{i-1,j}^\beta}{\Delta x} \\ \frac{\phi_{i,j+1}^\beta - \phi_{i,j-1}^\beta}{\Delta x} \end{pmatrix} \quad (6.87)$$

corresponding to the default discretization in the **Pace3D**-framework.

Alternatively, one can also extend the more face-centered discretization in Equation (6.82) in a similar manner by instead using the already calculated approximations of $\nabla\phi$ and $\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}$ on the cell phases and either setting

$$\begin{aligned} \left(\epsilon \frac{\partial a}{\partial\phi^\alpha} \right)_{i,j} &\approx 2\epsilon \sum_{\beta \neq \alpha} \frac{1}{4} \left[\left(\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i-\frac{1}{2},j}^{\alpha\beta}) \right) \cdot (\nabla\phi^\beta)_{i-\frac{1}{2},j} + \left(\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i+\frac{1}{2},j}^{\alpha\beta}) \right) \cdot (\nabla\phi^\beta)_{i+\frac{1}{2},j} \right. \\ &\quad \left. + \left(\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i,j-\frac{1}{2}}^{\alpha\beta}) \right) \cdot (\nabla\phi^\beta)_{i,j-\frac{1}{2}} + \left(\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}(\mathbf{q}_{i,j+\frac{1}{2}}^{\alpha\beta}) \right) \cdot (\nabla\phi^\beta)_{i,j+\frac{1}{2}} \right], \end{aligned} \quad (6.88)$$

or, in a somewhat cheaper manner which in addition reduces to the previous one if one is in fact isotropic, by setting

$$\begin{aligned} \left(\epsilon \frac{\partial a}{\partial\phi^\alpha} \right)_{i,j} &\approx 2\epsilon \sum_{\beta \neq \alpha} \frac{1}{2} \left[\left(\frac{\partial A^{\alpha\beta}}{\partial q_x^{\alpha\beta}}(\mathbf{q}_{i-\frac{1}{2},j}^{\alpha\beta}) \right) \frac{\phi_{i,j}^\beta - \phi_{i-1,j}^\beta}{\Delta x} + \left(\frac{\partial A^{\alpha\beta}}{\partial q_x^{\alpha\beta}}(\mathbf{q}_{i+\frac{1}{2},j}^{\alpha\beta}) \right) \frac{\phi_{i+1,j}^\beta - \phi_{i,j}^\beta}{\Delta x} \right. \\ &\quad \left. + \left(\frac{\partial A^{\alpha\beta}}{\partial q_y^{\alpha\beta}}(\mathbf{q}_{i,j-\frac{1}{2}}^{\alpha\beta}) \right) \frac{\phi_{i,j}^\beta - \phi_{i,j-1}^\beta}{\Delta y} + \left(\frac{\partial A^{\alpha\beta}}{\partial q_y^{\alpha\beta}}(\mathbf{q}_{i,j+\frac{1}{2}}^{\alpha\beta}) \right) \frac{\phi_{i,j+1}^\beta - \phi_{i,j}^\beta}{\Delta y} \right]. \end{aligned} \quad (6.89)$$

Remark 77. Which discretization of $\frac{\partial a}{\partial\phi^\alpha}$ is more favorable again depends on a variety of considerations. The one in Equation (6.87) is a priori the simplest, but has the potential disadvantage of requiring an additional calculation of the $\frac{\partial A^{\alpha\beta}}{\partial \mathbf{q}^{\alpha\beta}}$ in the cell-centers. In contrast, the discretizations in Equation (6.88) and (6.89) are again solely based on quantities already available on the faces. In addition, if combined with a buffering scheme, the contributions by each of the phases need only be calculated once as they are shared by the neighboring cells. This nevertheless has the disadvantage of introducing a spatial interdependence and some additional overhead through the buffering.

It should be stressed though that **none** of the formulations above is compatible with a discrete variational principle unless the contributions of the derivatives in the lateral direction in Equation (6.89) happen to drop out⁵³. Even though there is no principle difficulty in first defining

⁵³This is obviously the case for isotropic problems, but can also happen for some simple anisotropies such as e.g. an elliptic one described (in two dimensions) by

$$A^{\alpha\beta}(\mathbf{q}^{\alpha\beta}) = c_x^2 (q_x^{\alpha\beta})^2 + c_y^2 (q_y^{\alpha\beta})^2.$$

an “appropriate” discrete approximation of the energy \mathcal{F}_ϵ and then deriving the corresponding difference scheme through a differentiation with respect to the ϕ_i^α , this has a significant practical disadvantage as soon as this energy involves a gradient contribution calculated based on an averaged gradient, i.e. with a stencil spanning more than one grid spacing as this automatically entails a broader stencil.

This is most easily illustrated in an isotropic one-dimensional setting. Whereas the differentiation of an energy contribution of the form $\sum_j \left(\frac{\phi_{j+1} - \phi_j}{\Delta x} \right)^2$ with respect to ϕ_i leads to contributions by the two terms $j = i - 1$ and $j = i$ through $\frac{\phi_i - \phi_{i-1}}{\Delta x} \frac{1}{\Delta x}$ and $\frac{\phi_{i+1} - \phi_i}{\Delta x} \frac{-1}{\Delta x}$ and thus the standard second-order three-point stencil $-\frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2}$ of the negative Laplacian, an energy contribution of the form $\sum_j \frac{\phi_{j+1} - \phi_{j-1}}{2\Delta x}$ has the two contributions $\frac{\phi_i - \phi_{i-2}}{2\Delta x} \frac{1}{2\Delta x}$ and $\frac{\phi_{i+2} - \phi_i}{2\Delta x} \frac{-1}{2\Delta x}$ and thus instead leads to an undesirable five-point discretization $-\frac{\phi_{i+2} - 2\phi_i + \phi_{i-2}}{(2\Delta x)^2}$ of the same operator.

For the same reason, an energetic formulation based, as in the isotropic case, on energetic contributions by the phases (but now dependent on the full $\mathbf{q}^{\alpha\beta}$ -vectors there) would lead to the much broader stencil shown in the left of Figure 6.7, which in addition is not based on the very convenient structure of the discrete divergence operator as the difference of the fluxes through the four (resp. six in three space dimensions) cell faces. The other two more obvious choices are either the use of a purely cell-centered expression for the energy (leading to the slightly smaller stencil in the right of Figure 6.7), or the use of a “corner-centered” scheme based purely on short differences, which would allow to maintain the more narrow stencil in Figure 6.6. Both of the these latter choices are well-known to reduce to inherently unstable schemes in the isotropic case though, and are therefore likely not to be recommended.

Besides the higher computational complexity generally associated with broader stencils, they raise, in combination with the obstacle potential, an additional difficulty related with the discussion in Subsection 6.2.3. As already discussed there, by the discontinuity of the second derivative of the basic one-dimensional phasefield profile one expects a $\mathcal{O}(1)$ error in the discretization of the second derivative for any stencil crossing the transition from the interface to the bulk. While this was shown not to reduce the second-order convergence in the one-dimensional case, the degree to which this effect enters the calculation of the profile was nevertheless also seen to have a notable impact on the numerical precision. Similar effects will also arise in a higher-dimensional setting and will, with a broader stencil, affect more points depending on their relative positioning with respect to this transition⁵⁴.

Unless an energy-based discretization is necessary due to e.g. the use of a particular algorithm, a discretization of the gradient energy contributions as in Equation (6.86) and either Equation (6.87) or Equation (6.89) therefore seems to be the preferable choice.

◇

⁵⁴Recall that the numerically most benign situation in Section 6.2.3 was the one where the stencil of the last “mobile” point with $0 < \phi_i < 1$ is essentially completely within the expected interface region, whereas the - quite high - error at the first bulk-point is only seen indirectly in the error in the multipliers for the constraint. For a five-point stencil in a given direction, the one for the last inner point will always have to stretch across the transition into the bulk by roughly one grid spacing, and is therefore likely to be associated with significant numerical errors.

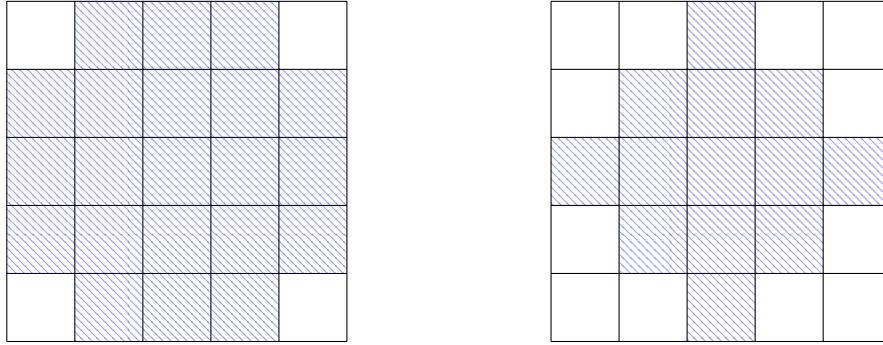


Figure 6.7: Difference stencils for a face-based (left) and cell-based (right) variational difference scheme for anisotropic gradient energy densities.

Dealing with Many Phases and the LROP Approach

In contrast to the mostly superficial difficulties from Subsection 6.3.3 (which are relatively easy to handle with in both a reduced and non-reduced fashion), a very crucial difficulty in the computational treatment of multiphase as compared to two-phase problems lies in the a priori very high computational and memory requirements associated with the large number of phases which are in many cases necessary for the simulation of realistic microstructures (with representative samples often consisting of thousands of grains).

Firstly, it is clear that, without further modifications, the memory requirements will increase linearly with the number N of phasefields required for the description of a given problem. Secondly, not only are there then also N equations to be solved (e.g. Equation (6.73) or Equation (6.75)), but they are in addition based on the calculation of two-phase interactions for a and even triple-phase interactions in the bulk potentials (see e.g. Equation (6.7) or (6.8)), and thus a priori lead to a cubic dependence of the calculation time on N .

A by now widely adapted manner for maintaining a roughly constant storage space and computational cost despite a large total number of phases is, at least in combination with an obstacle-type potential, the use of a **locally reduced order parameter (LROP)** approach ([40]). The underlying idea is that, based on stability considerations for multiphase-regions, one does not expect for regions with more than a few phases to be stable. If not initially present, such regions are therefore also not expected to arise naturally within the simulation, and one usually only has to deal with a very moderate maximal number l of phases which are actually present at any given point.

This is extremely useful from a computational point of view as it motivates replacing the storage of all phasefields with a scheme where one instead, for each computational cell, only stores the indices and the values of those phases with non-vanishing ϕ -values. While this a priori only allows reducing the memory requirements from $\mathcal{O}(N)$ to $\mathcal{O}(1)$ with respect to the number of phases, this storage scheme in addition also enables the reduction of the calculation time to a roughly constant one, regardless of the number of phases.

This is essentially due to the following two observations. On the one hand, the gradient energy contributions by construction vanish for all those phase-pairings for which at least one of the phases is locally constant at 0, as are the contributions to the equations of the other phases for all those α with $\phi^\alpha = 0$. Even though these phases themselves do experience a contribution by the other phases, the obstacle potential is designed to provide a strong counter-force to deviations of the local ϕ -values different from 0⁵⁵. This in combination with the general dominance of the

⁵⁵As well as from 1, but this is only of interest for at most one phase.

a - and w -terms in the phasefield functional is almost always (see Remark 79) sufficient to ensure that a locally “non-existent” phase which does not interact with its neighborhood through the non-local a -term will have to remain at 0. In addition, the multipliers due to the sum-constraint and the positivity constraint are such that the other phases remain unaffected by simply a priori skipping over the calculation of the new values of these phases as well as their contributions to the remaining equations. This can be made use of in terms of a classification scheme which identifies the phases actually present within a small neighborhood (depending upon the stencil) of a given cell and completely skips all calculations for the other ones. This classification is highly efficient in combination with the storage scheme above, as this simply amounts to running through the stored ϕ -indices of a fixed number of neighboring cells for checking which phases need to be considered, and thus is an operation of runtime $\mathcal{O}(1)$ with respect to the phases, with a reduction of the cost of the actual computations to the same order.

Remark 78. This preclassification can - and is, at least in the sense of which cells require any calculation at all, often worth the effort even in the two-phase case - in principle also be done even if all phases are stored. While this does allow a reduction of the actual calculation time to the same level as with the LROP-based scheme above, there are still two major drawbacks as compared to the storage scheme above once the number of phases grows even moderately large. Firstly, determining which phases should actually be considered as being present is still an operation of complexity $\mathcal{O}(N)$ as one is forced to run through all possible ones. Secondly, even if this classification does not require any real calculations (essentially reducing to a comparison of ϕ -values with 0), it can still incur a significant cost in run-time due to both a large number of conditional statements and the comparatively slow access to main memory⁵⁶. \diamond

While the storage could in principle also be done in an adaptive fashion with the scheme adjusting according to the local requirements, having an a priori estimate of the maximal number of phases one expects to locally coexist - either based on l as above or simply by experience for a particular type of setting - allows both a simplification in the implementation and an increase in the efficiency by using a less dynamic memory layout designed to provide sufficient space for locally storing a fixed number l' of phases only. In combination with the preclassification outlined above, this leads to a very performant implementation capable of handling an essentially arbitrary total number of phases with a roughly constant computational cost (see [40] and [67] for a more detailed discussion on various aspects of this approach).

Remark 79. It is clear that the approach above is not rigorous in the sense that one is not guaranteed to actually remain below any local prefixed number $l' < N$ of phases at all times. Even in settings specifically designed to lead to an initial violation of this bound, simulations with e.g. $l' = 6$ or $l' = 8$ typically lead to practically identical simulation results after a potential (but very short) transient which quickly forces an elimination of the “excess” number of phases. One example of such a setting is starting an N -phase simulation with a simple sharp interface between two given phases. This initial jump discontinuity in the phases leads to a discrete $\delta'(x)$ -type contribution and thus a term of order $\mathcal{O}\left(\frac{1}{(\Delta x)^2}\right)$ to the cells neighboring the interfaces due to the second-order derivative in the divergence-term $-\nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi}\right)$. Unless ϵ is chosen unrealistically small, this term will initially dominate the contribution by the bulk-potential and will, at least for an explicit time-discretization, enforce the appearance of **all** N phases through the Lagrange-multiplier Λ in Equation (6.73) resp. more indirectly but based on the same effect through the mobility matrix \mathbf{M} in Equation (6.75).

This can e.g. be seen by considering a very simple isotropic one-dimensional setting for which there is a sharp transition between a bulk-phase for $\alpha = 1$ up to the cell i to a bulk-phase for $\alpha = 2$ starting at $i + 1$, where, for simplicity of notation, it will additionally be assumed that all $\gamma^{\alpha\beta}$ are equal to the common value γ .

⁵⁶This depends of course also very significantly upon the details of the particular implementation and simulation setup (i.e. for example favorable cache effects). Regardless of this, it is an issue which is strongly mitigated using the LROP approach.

Together with the initially sharp transition, the evaluation of the simplified expression for $q^{\alpha\beta}$ in Equation (6.84) leads to the only non-zero terms begin $q_{i+\frac{1}{2}}^{12} = -q_{i+\frac{1}{2}}^{21} = \frac{1}{\Delta x}$, with $q_{i+\frac{1}{2}}^{\alpha\beta} = 0$ otherwise, and $q_{j+\frac{1}{2}}^{\alpha\beta} = 0$ regardless of the phases for all $j \neq i$. For the divergence-term $\epsilon \nabla \cdot \left(\frac{\partial a}{\partial \nabla \phi^\alpha} \right) = -2\gamma\epsilon \sum_{\beta \neq \alpha} q_{i+\frac{1}{2}}^{\alpha\beta}$ discretized as in Equation (6.80), the only non-zero contributions arising through the divergence-part in the cell i are given by $-\frac{\gamma\epsilon}{(\Delta x)^2}$ for the phase $\alpha = 1$ and $\frac{\gamma\epsilon}{(\Delta x)^2}$ for the phase $\alpha = 2$ since both $\phi_{i+\frac{1}{2}}^1$ and $\phi_{i+\frac{1}{2}}^2$ equal $\frac{1}{2}$. From Equation (6.69), it is also obvious that one has $\frac{\partial w}{\partial \phi^\alpha}(\phi_i) = 0$ if $\alpha = 0$ and $\frac{\partial w}{\partial \phi^\alpha}(\phi_i) = \frac{16\gamma}{\pi^2}$ for $\alpha \neq 1$. The only remaining relevant term⁵⁷ is given by $\epsilon \frac{\partial a}{\partial \phi^\alpha} = 2\gamma\epsilon \sum_{\beta \neq \alpha} q^{\alpha\beta} \frac{d\phi^\beta}{dx}$. Discretizing this expression for example as in Equation (6.82), the sum on the lower face in cell i is 0 and the upper one equals $\frac{2\gamma\epsilon}{(\Delta x)^2}$ if $\alpha = 1$ or $\alpha = 2$ and 0 otherwise. Averaging over both faces and combining the result with the remaining expressions above, it follows that one has

$$\epsilon \left(\nabla \cdot \frac{\partial a}{\partial \nabla \phi^\alpha} - \frac{\partial a}{\partial \phi^\alpha} \right) - \frac{1}{\epsilon} \frac{\partial w}{\partial \phi^\alpha} \hat{=} \begin{cases} \left(-\frac{\gamma\epsilon}{(\Delta x)^2} - \frac{\gamma\epsilon}{(\Delta x)^2} \right) - 0 = -\frac{2\gamma\epsilon}{(\Delta x)^2} & , \alpha = 1, \\ \left(\frac{\gamma\epsilon}{(\Delta x)^2} - \frac{\gamma\epsilon}{(\Delta x)^2} \right) - \frac{16\gamma}{\pi^2\epsilon} = -\frac{16\gamma}{\pi^2\epsilon} & , \alpha = 2, \\ -\frac{16\gamma}{\pi^2\epsilon} & \text{else} \end{cases}$$

and thus, if none of the μ^α differs from zero (as will be seen to be the case below)

$$\Lambda = -\frac{2\gamma\epsilon}{N(\Delta x)^2} - \frac{N-1}{N} \frac{16\gamma}{\pi^2\epsilon}.$$

In total, the right-hand sides in Equation (6.73) are therefore given by

$$-g^\alpha - \Lambda = \begin{cases} \frac{N-1}{N} \left(-\frac{2\gamma\epsilon}{(\Delta x)^2} + \frac{16\gamma}{\pi^2\epsilon} \right) & , \alpha = 1, \\ \frac{1}{N} \left(\frac{2\gamma\epsilon}{(\Delta x)^2} - \frac{16\gamma}{\pi^2\epsilon} \right) & , \alpha \neq 1, \end{cases}$$

and are strictly positive and the same for all phases with $\phi^\alpha = 0$, leading to the appearance of all phases. \diamond

Remark 80. Even though the example in the previous remark shows that there are certain risks associated with a priori assuming a result instead of performing the actual calculation, its intent is not at all to argue against the use of such an approach. Firstly, it is from a computational point of view essentially the only way one can reasonably perform simulations involving more than a few phases. As such, the use of a (somewhat more conservative) LROP-type scheme may even be justified when using a well-potential, despite the fact that one expects all phases to spread over the entire domain and thus the introduction of a persistent error due to the necessity of forcing phases with $\phi^\alpha \approx 0$ to zero. With respect to the other sources of error enforced by computational constraints - both numerical and in particular through the use of an often artificially large interface - this cut-off may in fact still lead to better results by allowing the use of higher resolutions and thus smaller interface widths.

Secondly, despite the algorithm delivering “wrong” results during one or a few time-steps, the unexpected appearance of all phases is essentially an unphysical artefact due to an interplay of the difficulties associated with modeling multiphase interactions combined with a very irregular initial setting and is therefore - even though strictly speaking incorrect - preferable from a physical point of view.

Nevertheless, one has to be very careful when skipping calculations or a priori excluding phases, as this can, depending on the setting, also either hide less well-known difficulties of the model or even lead to critical failures of the calculations⁵⁸. \diamond

⁵⁷Unless the driving force is very strong, it plays no role in the argument here.

⁵⁸Such failures are normally not due to a preclassification and the use of a restricted storage scheme alone, but its combination with a second not fully justified “optimization”.

Projection-Based Algorithms

A popular approach for the time-discretization of dynamic problems subject to additional constraints is the use of projection-based schemes⁵⁹. In these, one first obtains a prediction such as it would result from the equation obtained by partially or completely “ignoring” the constraint. Since this prediction need not be compatible with the constraints, it is then, in a second step, projected back onto the admissible set.

The basic idea of this two-step procedure is to separate the total problem into a sequence of two a priori simpler problems. The first one (i.e. obtaining the prediction) is then dependent upon the equation itself, whereas the second projection step is primarily tied to the structure of the constraint set and can therefore be performed in a relatively “generic” fashion without knowing the details of the various contributions arising in the equation itself.

While this cursory description is for the most part true in the sense that the projection operation does not have to take the precise nature of any fixed contribution to the equation (such as e.g. the actual expressions for the $g^\beta = \frac{d\mathcal{F}_\epsilon}{d\phi}$ if discretized in an explicit manner) into account, the combined updates must remain compatible with the underlying first-order necessary conditions characterizing the minimizer and thus, in the phasefield case, the KKT-condition in Equation (6.63). This implies that the projection must be constructed in a manner compatible with the way the prescribed forcings affect the prediction in the absence of the constraints. For this reason, while essentially independent of the “energy” underlying the gradients, it does depend crucially upon the chosen dynamics for the evolution.

This somewhat abstract difficulty as well as some related issues are best understood by considering an example. An in practice very relevant one in the phasefield context is given by an explicit time-discretization of the different choices of dynamics considered in Subsection 6.3.2. For the two choices of either a scalar proportionality of the dynamics as in Equation (6.72) or the mobility-matrix based one in Equation (6.77), an “explicit” discretization leads to

$$\tau^{(n)}_\epsilon \frac{\phi^{(n+1)} - \phi^{(n)}}{\Delta t} = \mathbf{r}^{(n)} - \Lambda^{(n')} \mathbf{e} + \boldsymbol{\mu}^{(n')} \quad (6.90)$$

resp.

$$\epsilon \frac{\phi^{(n+1)} - \phi^{(n)}}{\Delta t} = \mathbf{M}(\mathbf{r}^{(n)} + \boldsymbol{\mu}^{(n')}) \quad (6.91)$$

where $\tau^{(n)} := \tau(\phi^{(n)})$ and $\mathbf{r}^{(n)}$ denotes an appropriate discrete version of the negative gradient contributions from $\frac{d\mathcal{F}_\epsilon}{d\phi}$ in Equation (6.74). The n' for the multipliers is used to indicate some additional freedom in the choice and interpretation of their corresponding “time-step”.

One of the most popular algorithms for a constrained gradient-descent is to actually not consider (6.90) and (6.91) and the associated question of the choice of multipliers directly. Instead, the **projected gradient descent** algorithm corresponds to the the splitting approach above in its purest form, where one first generates a prediction $\hat{\phi}^{(n+1)}$ by completely ignoring the constraints through $\hat{\phi}^{(n+1)} = \phi^{(n)} + \frac{\Delta t}{\tau^{(n)}_\epsilon} \mathbf{r}^{(n)}$ resp. $\hat{\phi}^{(n+1)} = \phi^{(n)} + \Delta t \mathbf{M} \mathbf{r}^{(n)}$, and then projects this prediction back onto the admissible set. In order to be compatible with Equation (6.90), a direct comparison shows that the change $\delta\phi := \phi^{(n+1)} - \hat{\phi}^{(n+1)}$ induced by the projection operation has to satisfy

$$\phi^{(n+1)} = \phi^{(n+1)} - \hat{\phi}^{(n+1)} + \hat{\phi}^{(n+1)} = \delta\phi + \phi^{(n)} + \frac{\Delta t}{\tau^{(n)}_\epsilon} \mathbf{r}^{(n)} = \phi^{(n)} + \frac{\Delta t}{\tau^{(n)}_\epsilon} (\mathbf{r}^{(n)} - \Lambda^{(n')} \mathbf{e} + \boldsymbol{\mu}^{(n')})$$

⁵⁹Two particularly well-known examples for this approach are provided by the pressure-projection schemes in fluid-dynamical simulations and the popular projection-based approach in elasto-plastic simulations (where it is usually only used as a substep of a global more complex algorithm).

and thus $\delta\phi = \frac{\Delta t}{\tau^{(n)}\epsilon}(-\Lambda^{(n')}\mathbf{e} + \boldsymbol{\mu}^{(n')})$, whereas the analogous calculation for Equation (6.91) shows that $\delta\phi$ has to be of the form $\delta\phi = \mathbf{M}\boldsymbol{\mu}^{(n')}$.

This illustrates two points. On the one hand, the difference in structure of the increments due to the projection clearly indicates that the appropriate projection operation in both cases is different, and thus the dependence on the dynamics already discussed above. The appropriate projector operator for the (scalar) kinetic coefficient is the one based on the Euclidian norm (or some multiple thereof), since⁶⁰

$$\phi^{(n+1)} = \min_{\phi \in \mathcal{GS}^N} \frac{1}{2} |\phi - \hat{\phi}^{(n+1)}|^2 \Leftrightarrow \begin{cases} \phi^{(n+1)} = \hat{\phi}^{(n+1)} - \bar{\Lambda}\mathbf{e} + \bar{\boldsymbol{\mu}}, \\ \bar{\boldsymbol{\mu}} \geq 0, \bar{\boldsymbol{\mu}}^\alpha \phi^\alpha = 0, \\ \phi \in \mathcal{GS}^N, \end{cases}$$

and it thus leads to the desired structure of $\delta\phi$. In contrast, when the dynamics are based on the mobility-matrix \mathbf{M} , the projection should be chosen as the one with respect to the semi-norm induced by the pseudo-inverse \mathbf{M}^\dagger of \mathbf{M} , as, if $\hat{\phi}^{(n+1)}$ has a zero average, this implies that $\phi^{(n+1)} = \hat{\phi}^{(n+1)} + \mathbf{M}\bar{\boldsymbol{\mu}}$ and thus an increment $\delta\phi$ of the required form⁶¹.

On the other hand, it also explains one of the main reasons for the high popularity of the projected gradient descent algorithm, namely its particularly simple formal structure. Whereas a direct discretization of the dynamics in equations (6.72) or (6.77) as for example in equations (6.90) and (6.91) requires some considerations with regards to the multipliers, these can, in the explicit case, be completely disregarded and will be generated implicitly through the projection operator.

An alternative approach - consistent with a continuous interpretation if starting from a point satisfying the constraints - would be to instead choose Λ and $\boldsymbol{\mu}$ such that the evolution is restricted to the tangent cone of the admissible set at the current point, and thus, due to the scalar proportionality on the left-hand sides of equations (6.90) resp. (6.91), a projection of the right-hand sides onto this tangent space instead of the projection of the result onto the admissible set itself. This tangent space is given by all “directions” \mathbf{d} such that $\sum_\alpha d^\alpha = 0$ for satisfying the sum-constraint and $d^\alpha \geq 0$ if $\phi^\alpha = 0$ and arbitrary otherwise in order not to move a phase at 0 below zero.

Projecting \mathbf{r} with respect to the Euclidian norm, i.e. minimizing $\frac{1}{2}|\mathbf{d} - \mathbf{r}|^2$ over this admissible set, the resulting direction can similar to above be written as $\mathbf{d} = \mathbf{r} - \bar{\Lambda}\mathbf{e} + \bar{\boldsymbol{\mu}}$. Here $\bar{\boldsymbol{\mu}}$ is subject to two “complementarity” conditions, namely $\bar{\boldsymbol{\mu}}^\alpha = 0$ if $\phi^\alpha > 0$ - this is known in advance and is thus a purely artificial one for being able to use an (unnecessary) full vector of multipliers for simplifying the notation - and $\bar{\boldsymbol{\mu}}^\alpha d^\alpha = 0$ for all α with $\phi^\alpha = 0$. Similarly, a projection of $\mathbf{M}\mathbf{r}$ with the respect to the semi-norm induced by \mathbf{M}^\dagger leads to the analogous conclusion with $\mathbf{d} = \mathbf{M}(\mathbf{r} + \bar{\boldsymbol{\mu}})$ with $\bar{\boldsymbol{\mu}}$ subject to the same conditions.

Remark 81. This is in a sense the “most explicit” discretization possible as it amounts to a projection of the given right-hand sides \mathbf{r} and is thus a priori completely independent of $\phi^{(n+1)}$, and the multipliers will therefore be denoted with $n' = n$. Nevertheless, it has the disadvantage that, even though the projection of \mathbf{r} is a feasible direction, it may only be so for a very restricted time-step. More precisely, a negative total right-hand side for a phase α with $\phi^\alpha \approx 0$ can force this phase below zero in the next time-step, unless Δt is chosen smaller than the one defined by the **break-point** where $0 = (\phi^\alpha)^{(n)} + \frac{(\Delta t)_b^\alpha}{\tau^{(n)}\epsilon}(\mathbf{r}^{(n)} - \Lambda^{(n)} + \boldsymbol{\mu}^{(n)})$ and the scheme will therefore only lead to feasible values in the next time-step if $\Delta t \leq \min_\alpha (\Delta t)_b^\alpha$. Even though the evaluation of the right-hand side can thus be made completely independent of the next time-step, this simply

⁶⁰Due to the the convexity of the norm, the first-order condition on the right being both necessary and sufficient.

⁶¹In contrast to the Euclidian case, this statement is a little more technical. As will be discussed in more detail in Section 6.3.3, this projection is still well-defined despite the use of the semi-norm and the lack of actual invertibility of \mathbf{M} does in addition not cause any major practical difficulties.

amounts to shifting the responsibility for maintaining the feasibility to the choice of Δt . This strategy is not by itself a bad idea and is in a similar form used as part of some optimization algorithms such as e.g. the reduced gradient method [58]. The major disadvantage here though is that the choice of Δt has global effects as it will have to be the same for all cells - and thus also has to be chosen based on the global minimum of the breakpoint of all phases in all cells - in order to be admissible for an actual time-discretization. If there is no restriction to map a physical evolution (i.e. in the current setting if every cell were allowed to use their own time-step), this can be a relatively convenient way of exploring a potential solution on the “faces” of the admissible set. \diamond

The two choices above are clearly not the only ones possible. In particular, one can e.g. choose to explicitly maintain an estimate Λ in the prediction step within the gradient projection scheme, i.e. to set $\hat{\phi}^{(n+1)} = \phi^{(n)} + \frac{\Delta t}{r^{(n)}\epsilon}(\mathbf{r}^{(n)} - \Lambda)$ and then project the result back onto the Gibbs-simplex. As the sum-constraint is always “active”, this is easily shown to have no effect on the final result, regardless how good or bad the choice of Λ .

Nevertheless, using e.g. the estimate $\Lambda = \frac{1}{N} \sum_{\alpha} r^{\alpha}$ corresponding to the sum-constraint alone can be useful in practice as it automatically covers the multiwell-case and can, in a somewhat modified form, also be helpful in combination with a preclassification scheme such as in the LROP approach. In particular, this Λ is always a lower bound for the (scaled) actual multiplier resulting from the projection step⁶². As will be discussed in more detail in the next section, together with the non-negativity of $\boldsymbol{\mu}$, this ensures that any ϕ -value satisfying $\phi^{\alpha} \leq 0$ **before** the projection operation will necessarily also satisfy $\phi^{\alpha} = 0$ **after** the projection.

The situation for $\boldsymbol{\mu}$ is somewhat more difficult. In fact, as long as an estimate of $\boldsymbol{\mu}$ is compatible with the complementarity condition and such that, if included in the predictor step, it does not affect the inequality constraints which are ultimately active, the precise values are again irrelevant and will be “complemented” by the projection operation in a way which leads to the same final results for $\phi^{(n+1)}$. It is intuitively clear that this is not at all a given, since one can e.g. choose some ϕ^{α} at zero and force it to an arbitrarily high value in the prediction step by simply using a sufficiently large estimate for μ^{α} .

An admissible choice - closely related to the preclassification schemes - is the use of Λ and $\boldsymbol{\mu}$ as they would result from the projection of the gradient itself discussed above. The obvious advantage of this choice is that all phases for which this is compatible with the complementarity between d^{α} and μ^{α} (and for those phases only!), the ϕ -values initially at 0 will remain there and thus need not be updated.

Remark 82. As Remark 79 shows though, the correctness of this approach depends crucially upon these complementarity conditions actually being satisfied, since this example leads to a prediction where all phases at 0 have a positive right-hand side. Interpreted in terms of a “prediction step”, the (counter)example in this remark can be interpreted as a partial projection of the gradient with respect to all inequality constraints for the ϕ^{α} at zero except for the one which is classified as active due to its appearance in the neighboring cell. Due to the positivity of the right-hand side for all these phases in the presence of the sum-constraint, this projection would require **negative** μ^{α} ’s to ensure that they remain at zero, i.e. $d^{\alpha} = 0$. Checking whether complementarity actually holds would require maintaining the most expensive calculations due to the necessity of disposing of the right-hand sides r^{α} , and thus completely eliminates any advantages in terms of the computational cost.

Similar to the discussion in Section 6.3.3, it is always, in one form or another, the ability to perform this type of speculation in “good conscience” due to the favorable properties of the obstacle-potential which leads to phasefield codes being able to handle very large problem sizes - both in terms of the domain sizes and in particular with respect to the number of unknowns per “cell” in terms of N - as compared to many a priori simpler nonlinear programming problems. \diamond

⁶²A similar observation was already made for the steady-state case in Equation (6.65).

Remark 83. Note that the main feature making the projected gradient algorithm attractive in the phasefield case is that it is based on a relatively simple projection with a purely local admissible set and can thus essentially be evaluated on-the-fly as one runs through the domain to update the phasefield values. This in particular eliminates the necessity for storing the associated multipliers Λ and $\boldsymbol{\mu}$ ⁶³.

In contrast, a straightforward extension of this two-step procedure above to situations where the dynamics are modified through e.g. either the use of some semi-implicit or implicit time-stepping scheme or a “lightweight” preconditioning of the basic gradient descent scheme would not enjoy the same favorable features. In the simplest case, this leads to a scheme where a **non-local** s.p.d. matrix \mathbf{A}^{-1} essentially replaces the action of \mathbf{M} on the right-hand sides. Due to this, the appropriate projection is then a non-local one with respect to the norm induced by \mathbf{A} , which is on the one hand a significantly more difficult projection operation and on the other hand will generally necessitate actually storing the multipliers required “internally” by the projection.

Therefore, other algorithms are generally more appropriate in this case. One of the earlier popular choices is given by a nonlinear extension of the standard SOR-smoother for linear equations through a **projected SOR (pSOR)** algorithm (see [17] and e.g. [21] for a discussion in the phasefield context). The basic idea is the use of a point-relaxation process, which consists in considering each cell at a time and determining $\phi^{(n+1)}$ as the projection of the value ϕ^* which would allow to solve the equations within this cell (or, in the nonlinear case, an approximation thereof using e.g. a single Newton-step) if the values of the phasefield in all other cells are taken as fixed. This procedure can then be further accelerated by, instead of projecting ϕ^* itself, projecting the overrelaxed value $(1 - \omega)\phi^{(n)} + \omega\phi^*$. Even though this procedure does again require the use of a weighted projection instead of the Euclidian one for consistency with the multipliers, it maintains the main advantage of the strict locality of the projections.

More recently, various alternative approaches have been considered (see e.g. [60], [12], [32], [30], [31] and [33]) which avoid the main drawback of the pSOR algorithm, namely that the strict locality of this approach, while responsible for its simplicity, can also be the primary limiting factor for its convergence rate. \diamond

Some Algorithmic Aspects of the Projection onto the Gibbs-Simplex

In order to complete the description of the projected gradient algorithm above, it still remains to clarify how the projection operation of the prediction $\hat{\phi}$ can actually be performed.

This is straightforward when the only constraint is the sum-constraint $\sum_{\alpha=1}^N \phi^\alpha \stackrel{!}{=} 1$ as it suffices to subtract the average deviation from $\hat{\phi}$, i.e. $\mathcal{P}_{\Sigma_1}(\hat{\phi}) = \hat{\phi} - \frac{1}{N} \sum_{\alpha=1}^N \hat{\phi}^\alpha$.

Alternatively, at least if one has $\sum_{\alpha} (\phi^\alpha)_i^{(n)} = 1$, the same result can also be obtained by directly subtracting the average of the “right-hand sides” in the simplified form of the update rule (6.90) without the additional factor $\boldsymbol{\mu}$, i.e. by setting

$$\tau^{(n)} \epsilon \frac{\phi^{(n+1)} - \phi^{(n)}}{\Delta t} = \mathbf{r}^{(n)} - \left(\frac{1}{N} \left(\sum_{\alpha} (r^\alpha)^{(n)} \right) - 1 \right) \mathbf{e}$$

which then actually corresponds to the use of the projection of the gradient.

In contrast, when the predicted ϕ -values need to be projected back onto

$$\mathcal{GS} = \left\{ \boldsymbol{\phi} : \sum_{\alpha} \phi^\alpha = 1, 0 \leq \phi^\alpha, \alpha = 1, \dots, N \right\}, \quad (6.92)$$

the projection operation is generally more difficult and can in particular not be put into a convenient explicit formula.

⁶³While this is unlikely to be an issue in the case of the (phase-independent) multiplier Λ , $\boldsymbol{\mu}$ is a vector consisting of phase-specific entries and therefore potentially problematic when a large number of phases are present. As one of the principal advantages of explicit schemes is (besides their simplicity) their comparatively low memory requirement, this a quite convenient property.

The Euclidian Projector As already indicated above, the necessary and sufficient condition characterizing the Euclidian projection as the minimizer $\phi = \operatorname{argmin}_{\psi \in \mathcal{GS}} \frac{1}{2} |\psi - \hat{\phi}|^2$ are given by

$$\begin{cases} \phi = \hat{\phi} - \Lambda e + \mu, \\ \mu \geq 0, \mu^\alpha \phi^\alpha = 0, \\ \phi \in \mathcal{GS}^N. \end{cases} \quad (6.93)$$

A summation over the entries of the first equation together with $\sum_\alpha \phi^\alpha = 1$ implies in particular that the multiplier Λ for the sum-constraint is given by

$$\Lambda = \frac{1}{N} \sum_\alpha \left(\hat{\phi}^\alpha - 1 + \sum_\alpha \mu^\alpha \right), \quad (6.94)$$

and thus that Λ is **not** given by the simple formula corresponding to the sum-constraint alone unless the multipliers μ happens to vanish.

Nevertheless (a point which has also already been noted before) defining $\Lambda^{(0)} := \frac{1}{N} \sum_\alpha \hat{\phi}^\alpha - 1$, the positivity of the μ^α implies that one always has $\Lambda \geq \Lambda^{(0)}$. In addition, any phase α in the set of active constraints $\mathcal{A} := \{\alpha : \phi^\alpha = 0\}$, consisting of those which are actually at 0 after the projection, satisfies $0 = \hat{\phi}^\alpha - \Lambda + \mu^\alpha$ and thus

$$\mu^\alpha = \Lambda - \hat{\phi}^\alpha \geq 0. \quad (6.95)$$

In particular, any estimate $(\mu^\alpha)^{(n)} = \Lambda^{(n)} - \hat{\phi}^\alpha$ for μ^α based on Equation (6.95) and a value $\Lambda^{(n)}$ satisfying $\Lambda^{(n)} \leq \Lambda$ is never larger than the actual value of μ^α , i.e. $(\mu^\alpha)^{(n)} \leq \mu^\alpha$.

A repeated application of these two observations makes it possible to show that the following very simply algorithm will determine the correct projection:

Algorithm 1. Input: $\hat{\phi}$.

1. Set $\phi^{(0)} = \hat{\phi}$ and $\mathcal{A}^{(0)} = \emptyset$.
2. If $\phi^{(n)} \in \mathcal{GS}^N$, stop.
3. Else:
 - Find all new active phases $\delta\mathcal{A}^{(n)} = \{\alpha \notin \mathcal{A}^{(n)} : (\phi^\alpha)^{(n)} \leq 0\}$ and update the set of active constraints: $\mathcal{A}^{(n+1)} = \mathcal{A}^{(n)} \cup \delta\mathcal{A}^{(n)}$.
 - Calculate $\delta\Lambda^{(n)} = -\frac{1}{N - |\mathcal{A}^{(n+1)}|} \sum_{\alpha \in \delta\mathcal{A}^{(n)}} (\phi^\alpha)^{(n)}$.
 - Update ϕ : For all $\alpha \in \mathcal{A}^{(n+1)}$ set (resp. keep) $(\phi^\alpha)^{(n+1)} = 0$. For all inactive constraints, $\alpha \notin \mathcal{A}^{(n+1)}$, set $(\phi^\alpha)^{(n+1)} = (\phi^\alpha)^{(n)} - \delta\Lambda^{(n)}$.
 - Return to step 2.

One then has the following

Lemma 4. Given any $\hat{\phi}$, Algorithm 1 determines the vector $\phi = \operatorname{argmin}_{\psi \in \mathcal{GS}} \frac{1}{2} |\psi - \hat{\phi}|^2$ corresponding the Euclidian projection in at most $N - 1$ steps.

The sets $\mathcal{A}^{(n)}$ increase monotonically to the correct set \mathcal{A} of active constraints, and the values of the corresponding multipliers Λ and μ are given by

$$\Lambda = \sum_n \delta\Lambda^{(n)} = \frac{1}{N - |\mathcal{A}|} \left(\sum_{\alpha \notin \mathcal{A}} \hat{\phi}^\alpha - 1 \right) \quad \text{and} \quad \mu^\alpha = \begin{cases} 0 & \alpha \notin \mathcal{A}, \\ \Lambda - \hat{\phi}^\alpha & \alpha \in \mathcal{A}. \end{cases} \quad (6.96)$$

Proof. The crucial point here are some useful monotonicity-properties. Given any prediction $\mathcal{A}^{(n)}$ of the active set which does not contain a phase falsely classified as being zero - and this is certainly true for $\mathcal{A}^{(0)} = \emptyset$ - it follows from $\phi^\alpha = 0$ if $\alpha \in \mathcal{A}^{(n)}$ that the sum-constraint also has to be satisfied by the remaining $N - |\mathcal{A}^{(n)}|$ phases. A simple summation over all of these phases in the first line of Equation (6.93) shows that $1 = \sum_{\alpha \notin \mathcal{A}^{(n)}} \phi^\alpha = \sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - (N - |\mathcal{A}^{(n)}|)\Lambda + \sum_{\alpha \notin \mathcal{A}^{(n)}} \mu^\alpha$ and thus, similar to Equation (6.94), that

$$\Lambda^{(n)} = \frac{1}{N - |\mathcal{A}^{(n)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 + \sum_{\alpha \notin \mathcal{A}^{(n)}} \mu^\alpha \right). \quad (6.97)$$

As before, by the positivity of the μ^α , this implies that $\Lambda \geq \frac{1}{N - |\mathcal{A}^{(n)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 \right) =: \Lambda^{(n)}$. Setting $\phi^{(n+1)} := \hat{\phi}^{(n)} - \Lambda^{(n)} \mathbf{e}$, it follows that $\phi^{(n+1)} \geq \hat{\phi} - \Lambda \mathbf{e}$, where the inequality is to be understood componentwise. In particular, for all α with $(\phi^\alpha)^{(n+1)} \leq 0$ one necessarily also has $\hat{\phi}^\alpha - \Lambda \leq (\phi^\alpha)^{(n)} - \Lambda^{(n)} \leq 0$. This implies that the constraint of any such phase is actually active, i.e. $\alpha \in \mathcal{A}$ and $\phi^\alpha = 0$, since, assuming this is not the case, the complementary condition enforces $\mu^\alpha = 0$ and thus by the first-order necessary condition $\phi^\alpha = \hat{\phi}^\alpha - \Lambda$, contradicting the non-negativity of ϕ^α . Adding all such phases to the set of active constraints $\mathcal{A}^{(n+1)}$ at the next iteration therefore again leads to a set satisfying $\mathcal{A}^{(n+1)} \subset \mathcal{A}$.

This new estimate of the set of active constraints further leads to the new estimate

$$\Lambda^{(n+1)} = \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n+1)}} \hat{\phi}^\alpha - 1 \right) = \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 \right) - \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \in \delta \mathcal{A}^{(n)}} \hat{\phi}^\alpha \right)$$

for the Lagrange multiplier for the sum-constraint, where the summation was artificially split into those phases already in $\mathcal{A}^{(n)}$ and those in $\delta \mathcal{A}^{(n)} = \mathcal{A}^{(n+1)} \setminus \mathcal{A}^{(n)}$. Using

$$\begin{aligned} \frac{1}{N - |\mathcal{A}^{(n+1)}|} &= \frac{1}{N - |\mathcal{A}^{(n)}|} \frac{N - |\mathcal{A}^{(n)}|}{N - |\mathcal{A}^{(n)}| - |\delta \mathcal{A}^{(n)}|} = \frac{1}{N - |\mathcal{A}^{(n)}|} \left(1 + \frac{|\delta \mathcal{A}^{(n)}|}{N - |\mathcal{A}^{(n)}| - |\delta \mathcal{A}^{(n)}|} \right) \\ &= \frac{1}{N - |\mathcal{A}^{(n)}|} \left(1 + \frac{|\delta \mathcal{A}^{(n)}|}{N - |\mathcal{A}^{(n+1)}|} \right), \end{aligned}$$

the first term can be rewritten as

$$\frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 \right) = \underbrace{\frac{1}{N - |\mathcal{A}^{(n)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 \right)}_{=\Lambda^{(n)}} + \frac{|\delta \mathcal{A}^{(n)}|}{N - |\mathcal{A}^{(n+1)}|} \underbrace{\frac{1}{N - |\mathcal{A}^{(n)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 \right)}_{=\Lambda^{(n)}}$$

in terms of the previous estimate $\Lambda^{(n)}$, thus showing that

$$\begin{aligned} \Lambda^{(n+1)} &= \Lambda^{(n)} + \frac{|\delta \mathcal{A}^{(n)}|}{N - |\mathcal{A}^{(n+1)}|} \Lambda^{(n)} - \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \in \delta \mathcal{A}^{(n)}} \hat{\phi}^\alpha \right) = \Lambda^{(n)} - \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \in \delta \mathcal{A}^{(n)}} (\hat{\phi}^\alpha - \Lambda^{(n)}) \right) \\ &= \Lambda^{(n)} - \frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \in \delta \mathcal{A}^{(n)}} (\phi^\alpha)^{(n)} \right), \end{aligned}$$

where the last summand is precisely the definition $\delta \Lambda^{(n)} = -\frac{1}{N - |\mathcal{A}^{(n+1)}|} \left(\sum_{\alpha \in \delta \mathcal{A}^{(n)}} (\phi^\alpha)^{(n)} \right)$ in algorithm (1). As all $(\phi^\alpha)^{(n)}$ in $\delta \mathcal{A}^{(n)}$ are non-positive, it on the one hand is also clear that $\delta \Lambda^{(n)}$ is non-negative, and the $\Lambda^{(n)}$ indeed form an increasing sequence, and on the other hand, repeating this calculation, that $\Lambda^{(n)} = \sum_{m < n} \delta \Lambda^{(m)}$.

Finally, the algorithm is terminated once $\phi^{(n)}$ contains no further non-positive entries except for those in $\mathcal{A}^{(n)}$. In addition, the $\phi^{(n)}$ always satisfy the sum-constraint by the construction of $\delta\Lambda^{(n)}$ and one therefore has $\phi^{(n)} \in \mathcal{GS}^N$. All that remains to be verified is therefore that all the $(\mu^\alpha)^{(n)}$ defined as in Equation (6.96) based on $\Lambda^{(n)}$ and $\mathcal{A}^{(n)}$ are indeed positive. This is clear though by the monotonic increase of the $\Lambda^{(n)}$ and the fact that a phase is only added to $\mathcal{A}^{(m)}$ provided it satisfies, due to the equality $(\phi^\alpha)^{(m)} = \hat{\phi}^\alpha - \Lambda^{(m)}$ by the summation formula for the $\delta\Lambda$, $0 \geq \hat{\phi}^\alpha - \Lambda^{(m)} \geq \hat{\phi}^\alpha - \Lambda^{(n)}$ for all $n \geq m$, therefore ensuring $(\mu^\alpha)^{(n)} = \Lambda^{(n)} - \hat{\phi}^\alpha \geq 0$.

That this has to happen in at most n steps is obvious as, starting from $\mathcal{A}^{(0)} = \emptyset$, each iteration adds at least one phase with $\phi^\alpha = 0$, of which there can be at most $N - 1$ for any vector in \mathcal{GS}^N . \square

Remark 84. Note that, combining the description in Algorithm 1 and the proof of Lemma 4, there are in fact a variety of ways this algorithm can be implemented in an “equivalent” fashion. In particular, one can replace the incremental formulation in terms of the $\delta\Lambda^{(n)}$ and the $\phi^{(n)}$ by a direct evaluation based only on $\hat{\phi}$ and $\Lambda^{(n)} = \frac{1}{N - |\mathcal{A}^{(n)}|} \left(\sum_{\alpha \notin \mathcal{A}^{(n)}} \hat{\phi}^\alpha - 1 + \sum_{\alpha \in \mathcal{A}^{(n)}} \mu^\alpha \right)$ as in Equation (6.97). In fact, the $\phi^{(n)}$ are not really needed directly in the algorithm, since they only serve to judge whether $\hat{\phi}^\alpha - \Lambda^{(n)} \leq 0$ for any $\alpha \notin \mathcal{A}^{(n)}$ in order to determine the update $\delta\mathcal{A}^{(n)}$ of the set of active constraints. This can also be done directly based on this formula for $\Lambda^{(n)}$ and $\hat{\phi}$, with ϕ only being generated at the end through

$$\phi^\alpha = \begin{cases} 0 & \alpha \in \mathcal{A}^{(n)}, \\ \hat{\phi}^\alpha - \Lambda^{(n)} & \alpha \notin \mathcal{A}^{(n)}, \end{cases} \quad (6.98)$$

once there is no new phase added at the step n .

In addition, even though one can add all phases with $(\phi^\alpha)^{(n)} \leq 0$ to the new prediction of \mathcal{A} , one does not have to, and can e.g. only add the first such one that one encounters and then enter the next step of the algorithm (this being the original implementation in the **Pace3D**-framework), either in an incremental fashion using the $\delta\Lambda^{(n)}$ and the $\phi^{(n)}$ or in a non-incremental one based directly on $\hat{\phi}$ and $\Lambda^{(n)}$. Which version is more favorable computationally depends on a number of practical factors. While taking all newly active constraints into account at the same time does certainly deliver the lower number of total iterations, simply “skipping” the remaining checks and restarting with an updated (and increased) estimate of Λ increases the chances of more quickly identifying another phase which will ultimately have to be added to \mathcal{A} anyway and can thus also be advantageous.

Besides a pure iteration count and the number of arithmetic operations this involves, the actual run-time also depends on the time spend traversing the loops⁶⁴ in addition to the arithmetic operations themselves. \diamond

Remark 85. It should also be noted that the projection above (regardless of its actual implementation) interplays very favorably with a preclassification of the phases and an LROP-based approach. More precisely, if, instead of actually using all phases in this algorithm, one a priori can already exclude a number of those since they are already “known” to be zero, one can simply restrict the projection to the reduced subvector of phases for which this is not a priori clear without any effect on the result. In practical terms, this simply corresponds to taking only those phases into account for the algorithm, and replacing the total number N of phases with the number \tilde{N} of phases which are not a priori fixed at 0. \diamond

⁶⁴For example skipping over phases in $\mathcal{A}^{(n)}$ requires an additional conditional statement unless the phase is directly eliminated through an LROP-type indirect indexing. If based on this indirect indexing, eliminating a phase at 0 necessitates a rearrangement of indices, which can be an expensive operation if there are a large number of them and the one to eliminate appears very early.

The Weighted Projector Induced by the Mobility-Matrix M The projection is - both on a theoretical and practical level - somewhat more difficult when the dynamics are based on the use of the mobility-matrix M as in Equation (6.91). As discussed in the previous subsection, the compatibility of the projection-based approach (at least in the final steady state) - requires that the projection $\phi^{(n+1)}$ and the prediction $\hat{\phi}^{(n+1)}$ be related through $\phi^{(n+1)} = \hat{\phi}^{(n+1)} + M\mu^{(n+1)}$ with $\mu^{(n+1)}$ subject to the same complementarity conditions with respect to the (a priori unknown) vector $\phi^{(n+1)}$ as in the Euclidian case, at least if $\hat{\phi}$ is such that it satisfies the sum-constraint⁶⁵. In analogy with the case of invertible matrices⁶⁶ and dropping the time-indices corresponding to the “outer” algorithm, it will be seen below that this is essentially the first-order necessary condition corresponding to the minimization problem

$$\phi = \operatorname{argmin}_{\psi \in \mathcal{G}_{\mathcal{S}^N}} \frac{1}{2}(\psi - \hat{\phi}) \cdot M^\dagger(\psi - \hat{\phi}),$$

where M^\dagger is the pseudo-inverse of M .

Remark 86. Recall that $M^\dagger(\psi - \hat{\phi})$ for the Moore-Penrose pseudo-inverse is defined as the element ζ of smallest norm satisfying $M\zeta = \psi - \hat{\phi}$. Since the kernel of M consists of the constant vectors only as long as all $m^{\alpha\beta}$ are strictly larger than 0⁶⁷, this smallest norm-condition is not even necessary as long as $\hat{\phi}$ has the correct average as $\phi - \hat{\phi}$ is then a vector of zero average, meaning that the average of ζ plays no role for the semi-norm above. This is not true anymore if $\sum_\alpha \hat{\phi}^\alpha \neq 1$ though. \diamond

A straightforward differentiation of the corresponding Lagrangian $L(\phi, \Lambda, \mu) = \frac{1}{2}(\psi - \hat{\phi}) \cdot M^\dagger(\psi - \hat{\phi}) + \Lambda(\sum_\alpha \psi^\alpha - 1) - \mu \cdot \psi$ shows that the first-order necessary condition for this problem is given by $M^\dagger(\phi - \hat{\phi}) = -\Lambda e + \mu$ with satisfying the complementarity conditions $\mu \geq 0$, $\mu^\alpha \phi^\alpha = 0$. Multiplying by M , it follows that

$$MM^\dagger(\phi - \hat{\phi}) = M(-\Lambda e + \mu) = M\mu.$$

This is not quite the same as the $\phi = \hat{\phi} + M\mu$ as one has $MM^\dagger = \mathcal{P}_{\operatorname{Range}(M)}$ (see e.g. chapter 3 in [8]). By the fundamental theorem of linear algebra (Theorem 1) with $m = n$, $\operatorname{Range}(M) \oplus \operatorname{Ker}(M) = \mathbf{R}^n$ and since the kernel of M consists of the constant vectors only, it follows that $\phi - \hat{\phi}$ differs from $M\mu$ only through a constant vector of the form $-\Lambda e$, i.e. the the projected vector actually has to satisfy

$$\phi = \hat{\phi} + M\mu - \Lambda e. \quad (6.99)$$

Remark 87. It is not overly surprising that the sum constraint leads to the appearance of a multiplier Λ in a similar manner as in the Euclidian case, i.e. outside of M since M has no manner of changing the average which nevertheless has to be the correct one due to the admissible set.

⁶⁵While this is due to the structure of M normally being ensured by the dynamics, provided the initial phasefield-vectors are compatible with this constraint, this can cause some problems in combination with a too low number of phases which can be stored in the LROP-approach. Even if the storage space is sufficient for actually storing $\phi^{(n+1)}$, this need not be the case for $\hat{\phi}^{(n+1)}$ which may contain some phases below zero which would then be eliminated by the projection.

⁶⁶Recall that M is a singular matrix as it satisfies $M\mathbf{e} = \mathbf{0}$, i.e. it vanishes on the constant vectors.

⁶⁷This is a very natural assumption in the variational context as it is the only way of generally ensuring that ϕ becoming stationary in combination with μ as above implies that one has actually found a local minimizer. More precisely, if M only vanishes on the constant vectors, $M(-\mathbf{g} + \mu) = \mathbf{0}$, implies that $\mathbf{g} = -\Lambda e + \mu$ and thus the first-order necessary condition in Equation (6.63). Nevertheless, one seeming advantage of the mobility-based formulation in contrast to the one based on τ is that, in the prediction step, one can easily eliminate all interactions between certain phases by simply setting $m^{\alpha\beta} = 0$ without incurring a “division by zero”. Besides its problematic implications with respect to any potential minimizing or maximizing phasefield, this simplicity is highly misleading, since it causes some very tedious issues in the projection step. The strict positivity of the $m^{\alpha\beta}$, regardless of whether or not they are very small, will therefore be a standing assumption made here.

In contrast to the Euclidian case, the correct value of Λ can be determined upon entry into the projection operation, since - regardless of the values of $\boldsymbol{\mu}$ due to $\text{Range}(\mathbf{M}) = \text{Span}\{\mathbf{e}\}^\perp$ - a simple summation shows that $\Lambda = \sum_\alpha \hat{\phi}^\alpha - 1$. If $\hat{\phi}$ already satisfies the sum-constraint, it is then obvious that $\Lambda = 0$ and can thus safely be ignored. If not, it nevertheless suffices to replace $\hat{\phi}$ by $\hat{\phi} - \Lambda \mathbf{e}$ within the actual projection operation in order to obtain the correct result. \diamond

Remark 88. This is one of the points where vanishing values of $m^{\alpha\beta}$ can cause serious issues. Assuming e.g. there is a single phase β which is “disconnected” from the other ones, this adds an additional base vector \mathbf{e}_β to $\text{Ker}(\mathbf{M})$ and, by symmetry, to $\text{Range}(\mathbf{M})^\perp$. The conclusion above therefore needs to be modified to $\phi = \hat{\phi} + \mathbf{M}\boldsymbol{\mu} - \Lambda \mathbf{e} + \lambda^\beta \mathbf{e}_\beta$. As $\mathbf{e}_\beta \in \text{Ker}(\mathbf{M})^\perp$, there is no influence of $\boldsymbol{\mu}$ on the β -th entry of this equation, i.e. one has $\phi^\beta = \hat{\phi}^\beta - \Lambda + \eta^\beta$. This in itself essentially leaves λ^β undefined, provided it is such that ϕ^β satisfies the non-negativity constraint. It is only in combination with the definition of the pseudo-inverse - enforcing that $\mathbf{M}^\dagger(\phi - \hat{\phi})$ be the element of smallest norm satisfying this equation - that one is indirectly able to fix both Λ and η^β .

If the “basic” prediction of Λ is such that ϕ^β remains non-negative, the smallest deviation in norm is obtained for $\eta^\beta = 0$ as this is the optimal way of distributing a difference in the sum onto a vector (the norm involved in the pseudo-inverse being the standard Euclidian one here and not a weighted one) and the previous formula for Λ remains valid. If Λ moves ϕ^β out of the constraint set though - either directly due to $\phi^\beta < 0$ or indirectly through $\hat{\phi}^\beta - \Lambda > 1$ in combination with the sum-constraint and the constraints on the other phases - η^β has to be adjusted in a manner similar to the role of $\boldsymbol{\mu}$ in the Euclidian projection, except that one now has to keep two bounds in mind.

Even though this issue can in principle be handled and the phase β is afterwards unaffected by the determination of $\boldsymbol{\mu}$, this requires a rather complex “preprocessing” operation before considering the original problem of the determination of the $\boldsymbol{\mu}$. The situation is of course much worse when it is not a priori known that there is only a single disconnected phase but when there may be one or several groups of disconnected phases. \diamond

Whenever all mobilities are strictly positive, the projection can in fact be performed in a manner relatively similar to the Euclidian case. More precisely, it will be shown in Lemma 5 below that this can be achieved through

Algorithm 2. Input: $\hat{\phi}$.

1. Set $\phi^{(0)} = \hat{\phi}$ and $\mathcal{A}^{(0)} = \emptyset$.
2. If $\phi^{(n)} \in \mathcal{GS}^N$, stop.
3. Else:
 - Find all new active phases $\delta\mathcal{A}^{(n)} = \{\alpha \notin \mathcal{A}^{(n)} : (\phi^\alpha)^{(n)} \leq 0\}$ and update the set of active constraints: $\mathcal{A}^{(n+1)} = \mathcal{A}^{(n)} \cup \delta\mathcal{A}^{(n)}$.
 - Update the estimate for $\boldsymbol{\mu}$ by solving $\mathbf{M}_{\mathcal{A}^{(n+1)}\mathcal{A}^{(n+1)}}\boldsymbol{\mu}_{\mathcal{A}^{(n+1)}} = -\hat{\phi}_{\mathcal{A}^{(n+1)}}$, where $(\cdot)_{\mathcal{A}^{(n+1)}}$ denotes the restriction of the respective matrices and vectors to the index-set $\mathcal{A}^{(n+1)}$.
 - Update ϕ : For all $\alpha \in \mathcal{A}^{(n+1)}$ set (resp. keep) $(\phi^\alpha)^{(n+1)} = 0$. For all inactive constraints, set $(\phi^\alpha)^{(n+1)} = \hat{\phi}^\alpha + \sum_{\beta \in \mathcal{A}^{(n+1)}} M^{\alpha\beta} \mu^\beta = \hat{\phi}^\alpha - \sum_{\beta \in \mathcal{A}^{(n+1)}} m^{\alpha\beta} \mu^\beta$.
 - Return to step 2.

Remark 89. Note that, in contrast to Algorithm 1, it is now necessary to explicitly keep track of the values of the multipliers $\boldsymbol{\mu}$ as they affect the remaining ϕ -values differently due to the non-uniform weighting by the $m^{\alpha\beta}$. This is due to the fact that, unlike for the Euclidian case, the effect of each μ^α is not the same on all remaining phases due to the weighting and can therefore not be “summarized” through the action of a single multiplier Λ . \diamond

The following analogue of Lemma 4 holds:

Lemma 5. *Assume that the mobility matrix \mathbf{M} is of the form*

$$M^{\alpha\beta} = \begin{cases} \sum_{\gamma \neq \alpha} m^{\alpha\gamma} & \alpha = \beta, \\ -m^{\alpha\beta} & \text{else} \end{cases}$$

with $m^{\alpha\beta} = m^{\beta\alpha} > 0$ for all $\alpha \neq \beta$.

Then, given any $\hat{\phi}$ satisfying $\sum_{\alpha} \hat{\phi}^{\alpha} = 1$ or being being “replaced” by $\hat{\phi} - \frac{1}{N}(\sum_{\alpha=1}^N -1)e$, Algorithm 2 determines the weighted projection $\phi = \operatorname{argmin}_{\psi \in \mathcal{G}\mathcal{S}} \frac{1}{2}(\psi - \hat{\phi}) \cdot \mathbf{M}^{\dagger}(\psi - \hat{\phi})$ in at most $N - 1$ steps.

The sets $\mathcal{A}^{(n)}$ increase monotonically to the correct set \mathcal{A} of active constraints, each of the submatrices $\mathbf{M}_{\mathcal{A}^{(n)}\mathcal{A}^{(n)}}$ is invertible, and the corresponding multiplier μ is given by $\mu^{\alpha} = 0$ if $\alpha \notin \mathcal{A}$ and

$$\mu_{\mathcal{A}} = -\mathbf{M}_{\mathcal{A}\mathcal{A}}^{-1} \hat{\phi}_{\mathcal{A}} \quad (6.100)$$

otherwise, all entries of this vector are non-negative, and ϕ is given by $\phi = \hat{\phi} + \mathbf{M}\mu$.

Proof. The proof of this lemma is based on a very similar argument as the one of Lemma 4 in combination with the favorable structure of \mathbf{M} , namely verifying that the set \mathcal{A} is monotonically increasing without ever falsely classifying a constraint as active and that the vector μ in Equation (6.100) therefore has only non-negative entries and (by construction) is zero if $\phi^{\alpha} > 0$.

The crucial property of \mathbf{M} here is that it, even though \mathbf{M} itself is only a (singular) M -matrix, each of the submatrices $\mathbf{M}_{\mathcal{A}'\mathcal{A}'}$ with \mathcal{A}' strictly included in $\{1, \dots, N\}$ is a strictly diagonally dominant matrix with only non-positive off-diagonal entries. In fact, taking any such submatrix, the diagonal entry $M^{\alpha\alpha}$ is of the form $\sum_{\beta \neq \alpha} m^{\alpha\beta}$ with all $m^{\alpha\beta} > 0$, whereas the off-diagonal entries consist of the entries $-m^{\alpha\beta}$ for all $\beta \in \mathcal{A}'$. As \mathcal{A}' is assumed to be a strict subset of the set of all phase-indices, $M^{\alpha\alpha}$ is then clearly diagonally dominant as $m^{\alpha\beta} - \sum_{\mathcal{A}' \ni \beta \neq \alpha} |m^{\alpha\beta}|$ equals the sum over all $m^{\alpha\beta} \notin \mathcal{A}'$ which is strictly positive.

This implies firstly (see e.g. [55] and [76]) that $\mathcal{M}_{\mathcal{A}^{(n)}\mathcal{A}^{(n)}}$ is indeed invertible, and furthermore that all entries of $\mathcal{M}_{\mathcal{A}^{(n)}\mathcal{A}^{(n)}}^{-1}$ are non-negative. Even though this is by itself not sufficient to show that μ is non-negative as $-\hat{\phi}_{\mathcal{A}}$ may very well contain negative entries (i.e. “projected” phases which are initially above 0), it is sufficient to show the monotonous increase of the $\mu^{(n)}$ in Algorithm 2 and thus, as $\mu^{(0)}$ corresponding to $\mathcal{A}^{(0)} = \emptyset$ is all zero, that μ indeed only consists of non-negative entries.

Assuming therefore that $\mathcal{A}^{(n)}$ does not contain any phases falsely classified as active (and thus necessarily $|\mathcal{A}^{(n)}| < N$ by the sum-constraint) and that $\mu^{(n)} \geq 0$ - this trivially holding for $n = 0$ - it remains to show the same for the step $n + 1$. By the definition of the $\phi^{(n+1)}$ and $\mu^{(n+1)}$ in Algorithm 2 together with the invertibility of the submatrices, it follows that $\phi^{(n+1)}$ actually satisfies

$$\phi^{(n+1)} = \hat{\phi} + \mathbf{M}\mu^{(n+1)}. \quad (6.101)$$

On the one hand, since, with $\mu^{\alpha} = 0$ for $\alpha \notin \mathcal{A}^{(n+1)}$, the summation in the product $\mathbf{M}\mu^{(n+1)}$ automatically reduces to the summation over all phases in $\mathcal{A}^{(n+1)}$ and thus the formula for $(\phi^{\alpha})^{(n+1)}$ if $\alpha \notin \mathcal{A}^{(n+1)}$. On the other hand, $\mu^{(n+1)}$ is chosen precisely such that $\hat{\phi}_{\mathcal{A}^{(n+1)}} + \mathbf{M}_{\mathcal{A}^{(n+1)}\mathcal{A}^{(n+1)}}\mu^{(n+1)} = \mathbf{0}$ and thus implies that $\phi_{\mathcal{A}^{(n+1)}}^{(n+1)} = \mathbf{0}$ as imposed in the algorithm.

Since $\mu_{\{1, \dots, N\} \setminus \mathcal{A}^{(n)}}^{(n)} = \mathbf{0}$ by definition (and all relations are trivially true for $n = 0$), it follows by restricting Equation (6.101) to the set $\mathcal{A}^{(n+1)}$ that

$$\phi_{\mathcal{A}^{(n+1)}}^{(n)} = \hat{\phi}_{\mathcal{A}^{(n+1)}} + \mathbf{M}_{\mathcal{A}^{(n+1)}\mathcal{A}^{(n+1)}}\mu_{\mathcal{A}^{(n+1)}}^{(n)}$$

and thus, by subtracting this from Equation (6.101), that

$$(\phi^{(n+1)} - \phi^{(n)})_{\mathcal{A}^{(n+1)}} = \mathbf{M}_{\mathcal{A}^{(n+1)}\mathcal{A}^{(n+1)}}(\mu^{(n+1)} - \mu^{(n)})_{\mathcal{A}^{(n+1)}}.$$

As $\mathcal{A}^{(n)} \subset \mathcal{A}^{(n+1)}$, $\phi_{\mathcal{A}^{(n)}}^{(n+1)} = \phi_{\mathcal{A}^{(n)}}^{(n)} = \mathbf{0}$ (by the definition of the $(\phi^\alpha)^{(n)}$) and all entries in $\delta\mathcal{A}^{(n)} = \mathcal{A}^{(n+1)} \setminus \mathcal{A}^{(n)}$ satisfy $(\phi^{(n+1)} - \phi^{(n)})_{\delta\mathcal{A}^{(n)}} = \mathbf{0} - \phi_{\delta\mathcal{A}^{(n)}}^{(n)}$ with $\phi_{\delta\mathcal{A}^{(n)}}^{(n)} \leq \mathbf{0}$ (by the definition of $\delta\mathcal{A}^{(n)}$), it follows that the left-hand side consists of non-negative entries only. Together with the postivity of all entries of $\mathbf{M}_{\mathcal{A}^{(n+1)}\mathcal{A}^{(n+1)}}^{-1}$, it is then clear that the same also holds for $(\boldsymbol{\mu}^{(n+1)} - \boldsymbol{\mu}^{(n)})_{\mathcal{A}^{(n+1)}}$ and thus, with all other entries being zero, $\boldsymbol{\mu}^{(n+1)} - \boldsymbol{\mu}^{(n)}$.

Summarizing these arguments, the sets $\mathcal{A}^{(n)}$ increase by at least one entry per step unless there are no negative entries left in the prediction of the projection. All vectors satisfy the sum-constraint (potentially after a correction of their average as above) by the nature of \mathbf{M} and the $\boldsymbol{\mu}$ are monotonically increasing from $\mathbf{0}$ and therefore non-negative, showing that the algorithm converges to a solution of the necessary and here also sufficient condition in Equation (6.99). \square

Remark 90. It has to be observed that from a dynamic point of view, the mobility approach is in general **not** compatible with a preclassification scheme. Nevertheless, it can be shown that, once a steady state is reached and if the choice of skipping the calculations was correct (this being the crucial assumption of course), this state is compatible with the first-order necessary condition in Equation (6.63).

Even though this is therefore a potential problem only for the dynamics but not from a variational point of view, this point should not be underestimated as the motivation of the mobility-based formulation is precisely the modification of the dynamics, which, if used in an LROP-type fashion, gets mixed with a very “discrete” concept in terms of the preclassification.

What can in additionally generate some difficulties with the projection-based scheme here in combination with the LROP storage scheme is that, in contrast to the Euclidian projector, is that it is not sufficient for the calculation of the correct projection that a phase will be at 0 after the projection step. Instead, even if the result itself for this phase is known, it nevertheless necessary to also know the corresponding value in the prediction $\hat{\phi}$ as this difference has to be redistributed to the remaining phases in accordance with the weights in \mathbf{M} . Even if the storage is chosen sufficiently large to be able to store all “actually” occurring phases with non-zero values - i.e. those **after** the projection - this need not be the case before the projection as there may be a number of phases with negative values after the prediction step. Due to the discussion above, it is in this case not legitimate to simply truncate these to zero, since, even if this is their correct final value, the difference has to be known for the projection operation itself. \diamond

Chapter 7

Applications in the Material Sciences

The use of the phasefield method in the material sciences is based on combining its ability for capturing effects associated with surface energies with driving forces induced by an appropriate thermodynamic contribution to the phasefield functional together with evolution equations for the relevant additional fields such as e.g. the concentration or temperature. Within the **Pace3D**-framework, these thermodynamic contributions are generally described directly or indirectly in terms of a free energy density $f(\phi, \mathbf{c}, T, \dots)$, through which the underlying free energy densities for the bulk-phases are extended to the diffuse interface region.

Given that the surface and bulk free energies are basically fixed parameters based on measured quantities for the given materials, one of the key challenges in achieving accurate results based on the phasefield method therefore lies in the construction of a sufficiently accurate model of the material behavior **within** the transition region. Whereas the earlier models for the energy contributions within this region were usually based upon relatively simple interpolation procedures in terms of the available additional fields - such as the total concentration \mathbf{c} or, in the mechanical case, the total strains ϵ - the more recent modeling approach consists in including a larger degree of the underlying physics into the interpolation procedure.

As the chemically driven phase transformations where on the one hand historically the first setting within which this question was addressed in detail and on the other hand in many cases allows for obtaining a satisfactory model using a somewhat simpler modification than in the mechanical case, this setting will be considered first in Section 7.1. It will start with an outline of and the discussion of some practically relevant issues induced by a now well established approach for isothermal solidification problems developed, among others, in the works of [42], [25], [19] and [56]. Through a slight adaptation, this approach can in fact be extended to include the more general non-isothermal setting underlying the model in [52].

Section 7.2, will then discuss related modeling approaches within a mechanical setting. In contrast to the chemical setting, a simple translation of the ideas in the chemical case unfortunately does not result in equally satisfactory results due to a difference in the nature of the equilibrium conditions¹. Nevertheless, at least within a two-phase settings, a very satisfactory model based on the sharp interface mechanical jump conditions at interfaces has been developed in different forms in the works of [23], [51], [64] and [74]. Unlike in the chemical case, an extension to more than two phases involves a very fundamental difficulty closely related to the well-known singularities at singular points arising in the sharp interface setting.

¹This is somewhat misleading though as similar difficulties can in fact arise in the chemical setting as well depending on e.g. the boundary conditions.

Even though partial extensions to the multi-phase setting were proposed in, among others, [61], [62] and [74], these models loose, within multiphase regions, many of the favorable properties valid in the two-phase setting. The discussion of the multi-phasefield models together with some of these issues as well as modifcitations aiming to at least mitigate their impact will therefore be postponed until Section 7.2.4. Finally, Section 7.2.5 will outline a situation where both chemical and elastical effects arise in a coupled fashion, which will therefore be treated using a combination of both the chemical and mechanical models from the previous sections.

7.1 Quantitative Phasefield Models for Solidification

7.1.1 The ‘‘Traditional’’ Model

Before discussing some modifications through a more advanced modeling approach in the next section, it is instructive to recall the more traditional approach to alloy solidification problems as proposed in [52] and which parts of this model might be improved upon. As outlined in Section 3.2.1, it is based upon the entropy functional $\mathcal{S}_\epsilon(\phi, \mathbf{c}, e)$ from Equation (3.3) (repeated here for convenience)

$$\mathcal{S}_\epsilon(\phi, \mathbf{c}, e) = \int_{\Omega} s(\phi, \mathbf{c}, e) - (\epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi)) \, d\mathbf{x},$$

with the dependence of the thermodynamics of the material on the concentration \mathbf{c} and the local energy density e entering through the bulk entropy density $s(\phi, \mathbf{c}, e)$. $s(\phi, \mathbf{c}, e)$ can be related to the often more convenient free energy density $f(\phi, \mathbf{c}, T)$ through the Legendre-type transform pair

$$f(\phi, \mathbf{c}, T) = \inf_e \{e - Ts(\phi, \mathbf{c}, e)\} \quad \leftrightarrow \quad s(\phi, \mathbf{c}, e) = \inf_T \left\{ \frac{e}{T} - \frac{f(\phi, \mathbf{c}, T)}{T} \right\}$$

and it is thus sufficient to specify the bulk free energy density $f(\phi, \mathbf{c}, T)$ in order to specify $s(\phi, \mathbf{c}, e)$. Two examples for such a free energy density considered in [52] consist of an ideal-solution type model

$$f_{id}(\phi, \mathbf{c}, T) = \sum_{\alpha=1}^N \sum_{i=1}^K \left(c_i L_i^\alpha \frac{T - T_i^\alpha}{T_i^\alpha} h^\alpha(\phi) \right) + \sum_{i=1}^K \left(\frac{R}{v_m} T c_i \ln(c_i) \right) - c_v T (\ln(T) - 1) \quad (7.1)$$

and a subregular solution model, which, in its simplest form, reduces to the Redlich-Kister-type regular solution free energy density

$$f_r(\phi, \mathbf{c}, T) = f_{id}(\phi, \mathbf{c}, T) + \sum_{\alpha=1}^N \sum_{i=1}^K \sum_{j=1}^K c_i c_j A_{ij}^\alpha h^\alpha(\phi). \quad (7.2)$$

Here, L_i^α and T_i^α are the latent heats per unit volume and the melting temperatures respectively, R_g the gas constant, v_m the (constant) molar volume and c_v the specific heat per volume. In addition, the $\mathbf{A}^\alpha = (A_{ij}^\alpha)_{1 \leq i, j \leq K}$ in Equation (7.2) represent a set of binary interaction coefficients for each phase α .

In relation with the discussion above, the following observation can be made: Taking $\phi = \mathbf{e}^\alpha$, the free energy density within any bulk phase α is of the form

$$f_{id}^\alpha(\mathbf{c}, T) = c_i L_i^\alpha \frac{T - T_i^\alpha}{T_i^\alpha} + \frac{R}{v_m} T c_i \ln(c_i) - c_v^\alpha T (\ln(T) - 1),$$

and the above definition of the free energy density in Equation (7.1) therefore corresponds to the simple weighted average $f_{id}(\phi, \mathbf{c}, T) = \sum_{\alpha=1}^N f_{id}^\alpha(\mathbf{c}, T) h^\alpha(\phi)$ of the ones in the bulk-phases (a similar conclusion obviously also being valid for the regular solution case). In terms of the distinction between the phase-specific and average concentrations above, and, a priori, an analogous distinction between the phase-specific temperatures T^α and the average temperature T , this could also be interpreted as being the weighted average $f_{id}(\phi, \mathbf{c}, T) = \sum_{\alpha=1}^N f_{id}^\alpha(\mathbf{c}^\alpha, T^\alpha) h^\alpha(\phi)$ with the additional assumption that, within the interface region, one has $\mathbf{c}^\alpha = \mathbf{c}$ and $T^\alpha = T$ for all phases α .

Whereas this type of model was quite successful for purely temperature dependent problems, it was realized in e.g. the works of [42] and [56] that by replacing the simple assumption on the phase-specific concentrations by an energetically more favorable one, one could significantly

reduce model artefacts within the interface region. In fact, whereas the assumption of equal temperatures is compatible with the the global equilibrium condition corresponding to the equilibration of the temperature over the entire domain (and thus in particular the equilibrium of the temperatures T^α in each of the bulk-phases), this is not the case for the assumption on the concentration fields. Instead, for the entropy functional above, the equilibrium condition for the evolution of the (average) concentration field is **not** given by the equilibration of the concentrations. Instead, it is related to the chemical driving forces $\frac{\partial \mathcal{S}}{\partial \mathbf{c}} = -\frac{1}{T} \frac{\partial f}{\partial \mathbf{c}} = -\frac{\mu}{T}$ and thus, in combination with the equilibrium condition $T = \text{const}$ on the temperature, ultimately to the condition of the equilibration (in the appropriate sense²) of the chemical potentials.

This has motivated replacing the assumption of the equality of the phase-specific concentrations with a condition based upon the equality of the (reduced) chemical potentials. This condition in combination with phase-specific concentrations was first used explicitly in a two-phase setting by Kim, Kim and Suzuki in [42] and has subsequently been extended to the multiphase setting and reanalyzed in a number of papers, including [25], [56] and [19]. As first discussed in [25], this so-called called **local quasi-equilibrium** condition - which will be discussed in more detail in the next section - can in fact be interpreted as defining the free energy function in such a way that it minimizes the local free energy density $\sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha)h^\alpha(\phi)$ subject to a given phase-field vector ϕ and average concentration \mathbf{c} . Based upon this interpretation, it is quite natural that this model will behave more favorably with respect to the issue of excess interface energy contributions, thus explaining its large popularity despite its additional complexity.

Remark 91. With regards to the literature, there are two important points which should be stressed here.

Firstly, there is a subtle but crucial difference between the earlier work [41] and the subsequent seminal paper [42] by Kim, Kim and Suzuki. Whereas the latter may be considered as the motivation and basis for the subsequent extensions by the various authors mentioned above, the analysis in the former one is, despite its large formal similarity with the discussion in this section in terms of its reliance upon implicit functions, actually much more closely related to the analysis of a variational inconsistency (see Remark 97). Secondly, despite the recurrent appearance of the grand potential energy densities Ψ^α below (a point which has been amply discussed in the literature), the underlying variational principle - at least for the natural isolating boundary conditions on the temperature and concentration field (resp. the free energy in the isothermal case) - is still given by the maximization of the entropy (resp. minimization of the free energy). While the Ψ^α appear quite naturally in this model and there is indeed a close link between the total grand potential energy density Ψ and the free energy defined in equatin (7.3), these are not the same and should not be confused. \diamond

7.1.2 The “More Advanced” Free Energy Model in the Isothermal Case

Before discussing some practical implications and possible extensions in the following sections, the basic idea underlying the more advanced free energy model in the isothermal case will be outlined below. Even though the model can be considered to have originated in the work [42], the discussion will mostly follow the approach in [25], which, in the authors opinion, provides a much clearer interpretation of the resulting model.

This idea is in fact remarkably simple. As already indicated above, standard phasefield models for chemically driven solidification problems are based on using a single, “average” concentration field \mathbf{c} only. Within interface regions with several coexistent phases, there is a priori no reason to assume for the concentrations to be the same in all phases. It is therefore - in line with the

²There is a slight complication in the formulation of this equilibrium condition when based upon a model using the full concentration vector \mathbf{c} together with the sum-constraint $\sum_{i=1}^K c_i = 1$ instead of the more common formulation based upon a reduced concentration vector $\tilde{\mathbf{c}}$ consisting e.g. only out of the first $K - 1$ components, the last one then being, if required, recovered as a byproduct of the sum constraint. This basically technical point will be discussed in more detail below.

interpretation of the ϕ -values (resp. the values $h^\alpha(\phi)$) as the volume-fraction of a given phase - quite natural to interpret the total concentration as a weighted averages of the ones in the individual phases, i.e. to postulate $\mathbf{c} = \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi)$ and similarly for the total free energy density to be given by $f(\phi, \mathbf{c}, T) = \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi)$ as the result of the contributions by the phase-specific free energy densities f^α within the individual phases.

These hypotheses alone do not yet fully specify the model, as, for given values of the parameters ϕ , \mathbf{c} and T , there are many ways the total concentration could be distributed to the individual phases while still satisfying the constraint on their average. The simplest assumption to make - leading to models of the type in [52] - is to assume that $\mathbf{c}^\alpha = \mathbf{c}$ for all phases α . For a chemically isolated system, the underlying variational principle in terms of the minimization of the free energy functional $\mathcal{F}_\epsilon(\phi, \mathbf{c}, T)$ is the redistribution of the initially present total concentration $\int_\Omega \mathbf{c} d\mathbf{x}$ in the energetically most favorable manner. The simple but elegant idea underlying the more quantitative models is to reuse this same principle for determining the phase-specific concentrations in terms of the average one, i.e. to redistribute the (given) average concentration onto the individual phases in such a manner as to minimize the total local free energy density $\sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi)$.

Remark 92. From a physical point of view, it is clear that one would also like to, in addition to the condition of the concentration average, enforce the conditions $0 < c_i < 1$ and $0 < c_i^\alpha < 1$ on the concentration values. An important simplification as compared to the case of the phasefield equation arises from the fact that one can usually safely ignore the 0-1-bounds on the concentration values, i.e. given a total concentration satisfying $0 < c_i < 1$, one will automatically have $0 < c_i^\alpha < 1$ without having to enforce this condition explicitly. For example in the case of the ideal and regular solutions above, this will be ensured by the logarithmic contributions, whose derivatives tend to infinity as a concentration value approaches 0 or 1.

In other cases, such as e.g. when using a (significantly cheaper) quadratic approximation of the free energy contributions, this need not be the case anymore though. The same of course also holds when, due to e.g. an ill-suited time-stepping, the average concentration \mathbf{c} falls out of these bounds. As such cases are usually due to a failure of some other part of the approximation scheme, the box-constraints will (except for some hints in the footnotes) nevertheless be systematically neglected in the following as this, in contrast to the discussion in the phasefield case with obstacle potential, adds quite some complexity with a very limited benefit. \diamond

In a more explicit manner, the f -function is thus given as the solution of the parameterized minimization problem

$$f(\phi, \mathbf{c}, T) := \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c})} \left\{ \sum_{\alpha} f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\}, \quad (7.3)$$

where the admissible set $\mathcal{A}(\phi, \mathbf{c})$ is, for any vector \mathbf{c} satisfying $\sum_{i=1}^K c_i = 1$, given by³

$$\mathcal{A}(\phi, \mathbf{c}) = \left\{ (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c} \right\}. \quad (7.4)$$

Alternatively, one can integrate the sum-constraint $\sum_{i=1}^K c_i^\alpha = 1$ directly by eliminating one component (here e.g. the last one) as a function of the others. In terms of the **reduced**

³If the sum-constraint on the average concentration is not satisfied, the admissible set is obviously empty as it is then not possible to obtain \mathbf{c} as an average of phase-specific concentrations \mathbf{c}^α satisfying this constraint. If one in addition wants to maintain the box-constraints on the concentration values, this admissible set would have to be modified to

$$\mathcal{A}(\phi, \mathbf{c}) = \left\{ (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c}, \sum_{i=1}^K c_i^\alpha = 1 \forall \alpha, 0 \leq c_i^\alpha, i = 1, 2, \dots, K, \forall \alpha \right\},$$

which, as in the phasefield case, will automatically guarantee $c_i^\alpha \leq 1$ due to the sum-constraint on the c_i^α .

concentration vectors $\tilde{\mathbf{c}}$ and $(\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N}$ consisting of the (then independent) first $K - 1$ -components, i.e. $\tilde{\mathbf{c}} = (c_1, c_2, \dots, c_{K-1})$ and similarly for the \mathbf{c}^α , the minimization problem (7.3) can alternatively be expressed as

$$\tilde{f}(\phi, \tilde{\mathbf{c}}, T) := \min_{(\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N} \in \tilde{\mathcal{A}}(\phi, \tilde{\mathbf{c}})} \left\{ \sum_{\alpha} \tilde{f}^\alpha(\tilde{\mathbf{c}}^\alpha, T) h^\alpha(\phi) \right\}, \quad (7.5)$$

where $\tilde{f}^\alpha(\tilde{\mathbf{c}}^\alpha, T) := f^\alpha((\tilde{\mathbf{c}}^\alpha, 1 - \sum_{i=1}^{K-1} c_i^\alpha), T)$ and the admissible set $\tilde{\mathcal{A}}(\phi, \tilde{\mathbf{c}})$ is given by⁴

$$\tilde{\mathcal{A}}(\phi, \tilde{\mathbf{c}}) = \left\{ (\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) = \tilde{\mathbf{c}} \right\}. \quad (7.6)$$

An important preliminary observation for the free energy functions defined in Equation (7.3) resp. Equation (7.5) is given in the following

Lemma 6. *Provided the f^α are (strictly) convex functions of the phase-specific concentrations \mathbf{c}^α , and the $(h^\alpha(\phi))_{1 \leq \alpha \leq N}$ represent convex weighting coefficients (i.e. satisfy $0 \leq h^\alpha \leq 1$, $\sum_{\alpha=1}^N h^\alpha(\phi) = 1$), f is also a (strictly) convex function of \mathbf{c} .*

The same conclusion also holds with \tilde{f} and $\tilde{\mathbf{c}}$ replacing f and \mathbf{c} .

Proof. This can be seen by essentially the same argument as for the convexity of inf(imal)-convolutions (see e.g. [8]). Given a convex combination $\mathbf{c} = \lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2$, $0 \leq \lambda \leq 1$ of the concentration vectors \mathbf{c}_1 , \mathbf{c}_2 , any convex combination $\mathbf{c}^\alpha = \lambda \mathbf{c}_1^\alpha + (1 - \lambda) \mathbf{c}_2^\alpha$ with $(\mathbf{c}_1^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_1)$ and $(\mathbf{c}_2^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_2)$ satisfies

$$\sum_{\alpha} \mathbf{c}^\alpha h^\alpha(\phi) = \sum_{\alpha} (\lambda \mathbf{c}_1^\alpha + (1 - \lambda) \mathbf{c}_2^\alpha) h^\alpha(\phi) = \lambda \left(\sum_{\alpha} \mathbf{c}_1^\alpha h^\alpha(\phi) \right) + (1 - \lambda) \left(\sum_{\alpha} \mathbf{c}_2^\alpha h^\alpha(\phi) \right) = \lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2 = \mathbf{c},$$

$$\sum_{i=1}^K c_i^\alpha = \sum_{i=1}^K (\lambda (c_1)_i^\alpha + (1 - \lambda) (c_2)_i^\alpha) = \lambda \sum_{i=1}^K (c_1)_i^\alpha + (1 - \lambda) \sum_{i=1}^K (c_2)_i^\alpha = \lambda + (1 - \lambda) = 1$$

and $0 \leq c_i^\alpha$, and thus $\lambda \mathcal{A}(\phi, \mathbf{c}_1) + (1 - \lambda) \mathcal{A}(\phi, \mathbf{c}_2) \subset \mathcal{A}(\phi, \lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2)$. Using this set-inclusion and the assumed convexity of f^α in terms of the \mathbf{c}^α in the definition (7.3) shows that

$$\begin{aligned} f(\phi, \mathbf{c}, T) &= \min_{\{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2)\}} \left\{ \sum_{\alpha} f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\} \\ &\leq \min_{\{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \lambda \mathcal{A}(\phi, \mathbf{c}_1) + (1 - \lambda) \mathcal{A}(\phi, \mathbf{c}_2)\}} \left\{ \sum_{\alpha} f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\} \\ &= \min_{\{(\mathbf{c}_1^\alpha, \mathbf{c}_2^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_1) \times \mathcal{A}(\phi, \mathbf{c}_2)\}} \left\{ \sum_{\alpha} f^\alpha(\lambda \mathbf{c}_1^\alpha + (1 - \lambda) \mathbf{c}_2^\alpha, T) h^\alpha(\phi) \right\} \\ &\leq \min_{\{(\mathbf{c}_1^\alpha, \mathbf{c}_2^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_1) \times \mathcal{A}(\phi, \mathbf{c}_2)\}} \left\{ \sum_{\alpha} \left(\lambda f^\alpha(\mathbf{c}_1^\alpha, T) + (1 - \lambda) f^\alpha(\mathbf{c}_2^\alpha, T) \right) h^\alpha(\phi) \right\}. \end{aligned}$$

⁴Now, maintaining the box-constraints would require modifying $\tilde{\mathcal{A}}$ to

$$\tilde{\mathcal{A}}(\phi, \tilde{\mathbf{c}}) = \left\{ (\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) = \tilde{\mathbf{c}}, 0 \leq c_i^\alpha \leq 1, i = 1, 2, \dots, K - 1, \forall \alpha \right\},$$

where, this time, it is necessary to enforce the upper bound of 1 on the c_i^α as the sum-constraints $\sum_{i=1}^K c_i^\alpha = 1$ has disappeared.

As this last minimization problem is separable⁵, it can equivalently be replaced by

$$\lambda \underbrace{\min_{\{(\mathbf{c}_1^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_1)\}} \left\{ \sum_{\alpha} f^\alpha(\mathbf{c}_1^\alpha, T) h^\alpha(\phi) \right\}}_{=f(\phi, \mathbf{c}_1, T)} + (1-\lambda) \underbrace{\left\{ \min_{\{(\mathbf{c}_2^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}_2)\}} \left\{ \sum_{\alpha} f^\alpha(\mathbf{c}_2^\alpha, T) h^\alpha(\phi) \right\} \right\}}_{=f(\phi, \mathbf{c}_2, T)}$$

and thus $f(\phi, \mathbf{c}, T) \leq \lambda f(\phi, \mathbf{c}_1, T) + (1-\lambda)f(\phi, \mathbf{c}_2, T)$ and the mapping $\mathbf{c} \mapsto f(\phi, \mathbf{c}, T)$ is indeed convex. \square

As discussed in Subsection 5.1, the FONC for the minimizer in Equation (7.3) can e.g. be obtained by considering the Lagrangian⁶

$$L\left(\phi, \mathbf{c}, T, (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, \boldsymbol{\mu}, \hat{\boldsymbol{\lambda}}\right) = \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) - \boldsymbol{\mu} \cdot \left(\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) - \mathbf{c} \right) - \sum_{\alpha=1}^N \hat{\lambda}^\alpha \left(\sum_{i=1}^K c_i^\alpha - 1 \right)$$

with $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^N$. $\boldsymbol{\mu}$ could in principle be taken as an arbitrary vector in \mathbb{R}^K , but, since the sum-constraint $\sum_{\alpha=1}^N c_i^\alpha h^\alpha(\phi) = c_i$ is in fact automatically fulfilled for all K components as soon as it holds for $K-1$ components⁷, this would lead to an indeterminacy as $\boldsymbol{\mu}$ would be capable of “testing” a property which is already ensured through the multipliers $(\lambda^\alpha)_{1 \leq \alpha \leq N}$. In order to avoid this issue, it may seem preferable to choose $\boldsymbol{\mu}$ in a suitable $(K-1)$ -dimensional subspace of \mathbb{R}^K . One possibility here, closely related to the reduced formulation, is to choose $\boldsymbol{\mu}$ such that $\mu_K = 0$. A second, more “symmetric” choice is obtained by enforcing a zero average on $\boldsymbol{\mu}$, i.e. by requiring that $\sum_{i=1}^K \mu_i = 0$. Alternatively, and this is in many ways the most natural choice, one may also fix $\boldsymbol{\mu}$ “indirectly” by simply requiring that $\boldsymbol{\mu} = \sum_{\alpha=1}^N \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha} h^\alpha(\phi)$, which in particular enforces $\sum_{\alpha=1}^N \lambda^\alpha h^\alpha(\phi) = 0$.

Remark 93. Each choice clearly leads to a different, though closely related, set of multipliers $(\boldsymbol{\mu}, \lambda^\alpha)$ as will be discussed in Subsection 7.1.5 below. The advantage of the last choice is that it carries, in terms of the average values, the standard role of the sensitivity with respect to the sum-constraint on the total concentration. The first two choices on the other hand “shift” this responsibility to the \mathbf{c}^α , therefore emphasising the role of the λ^α on the effective equation. This has no influence on the actual result, since it does not matter how the average sum-constraint on \mathbf{c} is ultimately enforced, but favors the last description on “aesthetic” grounds. \diamond

⁵It consists in the minimization of two sums of functions depending on independent arguments in independent sets. There is therefore no “competition” between the two summands, and the minimizer is obviously achieved when each of them is minimal individually.

⁶Together with the box-constraint, this needs to be modified to

$$L\left(\phi, \mathbf{c}, T, (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, \boldsymbol{\mu}, \hat{\boldsymbol{\lambda}}, \boldsymbol{\zeta}\right) = \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) - \boldsymbol{\mu} \cdot \left(\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) - \mathbf{c} \right) - \sum_{\alpha=1}^N \hat{\lambda}^\alpha \left(\sum_{i=1}^K c_i^\alpha - 1 \right) - \sum_{\alpha=1}^N \sum_{i=1}^K \zeta_i^\alpha c_i^\alpha$$

with $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^N$ and $\boldsymbol{\zeta} \in \mathbb{R}_+^{N \times K}$.

⁷Assuming for example that it holds for the first $K-1$ components and as the concentrations are required to sum to one, one has

$$\sum_{\alpha=1}^N c_K^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^N \left(1 - \sum_{i=1}^{K-1} c_i^\alpha \right) h^\alpha(\phi) = \sum_{\alpha=1}^N h^\alpha(\phi) - \sum_{i=1}^{K-1} c_i^\alpha h^\alpha(\phi) = 1 - \sum_{i=1}^{K-1} c_i = c_K.$$

In contrast, the modified Lagrangian for the reduced problem (7.5) is given by⁸

$$\tilde{L}(\phi, \tilde{\mathbf{c}}, T, (\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N}, \tilde{\boldsymbol{\mu}}) = \sum_{\alpha=1}^N f^\alpha(\tilde{\mathbf{c}}^\alpha, T) h^\alpha(\phi) - \tilde{\boldsymbol{\mu}} \cdot \left(\sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) - \tilde{\mathbf{c}} \right) \quad (7.8)$$

where now $\tilde{\boldsymbol{\mu}} \in \mathbb{R}^{K-1}$, and the multiplier $\hat{\boldsymbol{\lambda}}$ for the phase-specific sum-constraints has vanished. As both definitions (7.3) and (7.5) only differ in the representation of the minimization problem, it is clear that the numerical values of $\tilde{f}(\phi, \tilde{\mathbf{c}}, T)$ and $f(\phi, \mathbf{c}, T)$ will actually agree for any admissible concentration, i.e. one has $f(\phi, (\tilde{\mathbf{c}}, 1 - \sum_{i=1}^{K-1} c_i), T) = \tilde{f}(\phi, \tilde{\mathbf{c}}, T)$. Nevertheless, the different notations for both problems will be maintained in order to simplify the distinction between the two approaches.

Differentiation of L w.r.t. to the phase-specific concentration vectors \mathbf{c}^α leads to the system of equations

$$h^\alpha(\phi) \left(\frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, T) - \boldsymbol{\mu} \right) - \hat{\lambda}^\alpha \mathbf{e} = h^\alpha(\phi) (\boldsymbol{\mu}^\alpha - \boldsymbol{\mu}) - \hat{\lambda}^\alpha = \mathbf{0},$$

together with the sum-constraint $\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c}$, where $\boldsymbol{\mu}^\alpha := \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, T)$ and $\mathbf{e} = (1, 1, \dots, 1)$ is a K -dimensional vector of ones. Similarly, with $\tilde{\boldsymbol{\mu}}^\alpha := \frac{\partial \tilde{f}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}(\tilde{\mathbf{c}}^\alpha, T)$, the FONC for a minimizer of the reduced function in Equation (7.5) is given by the set of equations

$$h^\alpha(\phi) \left(\frac{\partial \tilde{f}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}(\tilde{\mathbf{c}}^\alpha, T) - \tilde{\boldsymbol{\mu}} \right) = h^\alpha(\phi) (\tilde{\boldsymbol{\mu}}^\alpha - \tilde{\boldsymbol{\mu}}) = \mathbf{0}$$

together with $\sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) = \tilde{\mathbf{c}}$.

A first observation to be made is that in both cases the conditions on the derivatives of the phase-specific f -functions **only** pertain to those phases for which $h^\alpha(\phi) \neq 0$, whereas the \mathbf{c}^α (resp. $\tilde{\mathbf{c}}^\alpha$) for the phases with $h^\alpha(\phi) = 0$ can in principle be chosen arbitrarily within the respective admissible sets⁹ without affecting either the value of f or $\sum_{\alpha=1}^N h^\alpha(\phi) \mathbf{c}^\alpha$ (resp. the reduced versions thereof). In the non-reduced case, the trivial choice $\lambda^\alpha = 0$ is then clearly an admissible choice for the multiplier λ^α for all vanishing phases, such that it suffices to focus only on those phases with $h^\alpha(\phi) \neq 0$ ¹⁰.

Defining the set $\mathcal{P}_p := \{\alpha \in 1, 2, \dots, N : h^\alpha(\phi) \neq 0\}$ of phases which are actually “present” at a given point, and for convenience scaling $\hat{\lambda}^\alpha, \alpha \in \mathcal{P}_R$ by $h^\alpha(\phi) \neq 0$, it is therefore in practice

⁸Or, in the case with box-constraints, by

$$\begin{aligned} \tilde{L}(\phi, \tilde{\mathbf{c}}, T, (\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N}, \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\zeta}}^-, \tilde{\boldsymbol{\zeta}}^+) &= \sum_{\alpha=1}^N f^\alpha(\tilde{\mathbf{c}}^\alpha, T) h^\alpha(\phi) - \tilde{\boldsymbol{\mu}} \cdot \left(\sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) - \tilde{\mathbf{c}} \right) \\ &\quad - \sum_{\alpha=1}^N \sum_{i=1}^{K-1} (\tilde{\zeta}^-)_i^\alpha c_i^\alpha - \sum_{\alpha=1}^N \sum_{i=1}^{K-1} (\tilde{\zeta}^+)_i^\alpha (1 - c_i^\alpha) \end{aligned} \quad (7.7)$$

with $\tilde{\boldsymbol{\mu}} \in \mathbb{R}^{K-1}$, $\tilde{\boldsymbol{\zeta}}^\pm \in \mathbb{R}_+^{N \times (K-1)}$. The introduction of two sets of multipliers $\tilde{\boldsymbol{\zeta}}^\pm$ is necessary as the inequality $c_i^\alpha \leq 1$ is not automatically ensured by the combination of the positivity and sum constraint anymore.

⁹This will be discussed in a little more detail in Subsection 7.1.3 when considering the derivative w.r.t. ϕ .

¹⁰Here, one could again include the box-constraints by replacing the FONC in the non-reduced case to

$$h^\alpha(\phi) \left(\frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, T) - \boldsymbol{\mu} \right) - \hat{\lambda}^\alpha \mathbf{e} - \boldsymbol{\zeta}^\alpha = h^\alpha(\phi) (\boldsymbol{\mu}^\alpha - \boldsymbol{\mu}) - \hat{\lambda}^\alpha - \boldsymbol{\zeta}^\alpha = \mathbf{0},$$

which then, in addition to the sum constraint, needs to be complemented by the box-constraints as well as $\zeta_i^\alpha \geq 0$ and the complementarity condition $\zeta_i^\alpha c_i^\alpha = 0$. Similarly, the reduced formulation would lead to the condition

$$h^\alpha(\phi) \left(\frac{\partial \tilde{f}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}(\tilde{\mathbf{c}}^\alpha, T) - \tilde{\boldsymbol{\mu}} \right) - (\tilde{\zeta}^-)^\alpha + (\tilde{\zeta}^+)^\alpha = h^\alpha(\phi) (\tilde{\boldsymbol{\mu}}^\alpha - \tilde{\boldsymbol{\mu}}) - (\tilde{\zeta}^-)^\alpha + (\tilde{\zeta}^+)^\alpha = \mathbf{0}$$

and the additional restrictions $(\tilde{\zeta}^\pm)_i^\alpha \geq 0$ and $(\tilde{\zeta}^-)_i^\alpha c_i^\alpha = 0$, $(\tilde{\zeta}^+)_i^\alpha (1 - c_i^\alpha) = 0$.

For phases with $h^\alpha(\phi) = \mathbf{0}$, admissible choices for the multipliers would then be given by $(\lambda^\alpha, \boldsymbol{\zeta}^\alpha) = (\mathbf{0}, \mathbf{0})$ (resp. $\tilde{\boldsymbol{\zeta}}^- = \tilde{\boldsymbol{\zeta}}^+ = \mathbf{0}$).

sufficient to consider the system

$$\begin{cases} \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, T) - \boldsymbol{\mu} - \lambda^\alpha \mathbf{e} = \boldsymbol{\mu}^\alpha(\mathbf{c}^\alpha, T) - \boldsymbol{\mu} - \lambda^\alpha \mathbf{e} = \mathbf{0}, & \alpha \in \mathcal{P}_p, \\ \sum_{i=1}^K c_i^\alpha = 1, & \alpha \in \mathcal{P}_p, \\ \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\boldsymbol{\phi}) = \mathbf{c} \end{cases} \quad (7.9)$$

with $\lambda^\alpha = \frac{\tilde{\lambda}^\alpha}{h^\alpha(\boldsymbol{\phi})}$ in the non-reduced case and

$$\begin{cases} \frac{\partial \tilde{f}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}(\tilde{\mathbf{c}}^\alpha, T) - \tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}^\alpha(\tilde{\mathbf{c}}^\alpha, T) - \tilde{\boldsymbol{\mu}} = \mathbf{0}, & \alpha \in \mathcal{P}_p, \\ \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\boldsymbol{\phi}) = \tilde{\mathbf{c}} \end{cases} \quad (7.10)$$

in the reduced case.

The system (7.10) with the **reduced chemical potential** or **diffusion potential** $\tilde{\boldsymbol{\mu}}$ is (except for the temporary restriction to $\alpha \in \mathcal{P}_p$) precisely the same condition also used in the models by [42], [25], [56] and [19]. As

$$\frac{\partial \tilde{f}^\alpha(\tilde{\mathbf{c}}^\alpha, T)}{\partial \tilde{\mathbf{c}}^\alpha} = \frac{\partial f^\alpha(\tilde{\mathbf{c}}^\alpha, 1 - \sum_{i=1}^{K-1} c_i^\alpha, T)}{\partial \tilde{\mathbf{c}}^\alpha} = \left(\frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_i^\alpha} \right)_{1 \leq i \leq K-1} - \frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_K^\alpha} \tilde{\mathbf{e}}, \quad (7.11)$$

where $\tilde{\mathbf{e}} = (\mathbf{e})_{1 \leq i \leq K-1}$, it follows from equations (7.10) and (7.9) that $\tilde{\boldsymbol{\mu}}$ can be obtained from $(\boldsymbol{\mu}, \lambda^\alpha)$ - independent of the specific phase α and the particular $(K-1)$ -dimensional subspace chosen for $\boldsymbol{\mu}$ - through the relation

$$\begin{aligned} \tilde{\mu}_i = \tilde{\mu}_i^\alpha &= \frac{\partial \tilde{f}^\alpha(\tilde{\mathbf{c}}^\alpha, T)}{\partial c_i^\alpha} = \frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_i^\alpha} - \frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_K^\alpha} = \mu_i^\alpha - \mu_K^\alpha \\ &= (\mu_i + \lambda^\alpha) - (\mu_K + \lambda^\alpha) = \mu_i - \mu_K. \end{aligned} \quad (7.12)$$

Conversely, as $\frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_K^\alpha} = \mu_K + \lambda^\alpha$, the solution satisfies¹¹

$$\frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_i^\alpha} = \mu_i + \lambda^\alpha = \mu_i + \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_K^\alpha} - \mu_K = \tilde{\mu}_i + \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_K^\alpha}.$$

In the simplest case $\boldsymbol{\mu}$ is, as suggested above, chosen such that $\mu_K = 0$, in which case the previous equation implies that

$$\lambda^\alpha = \frac{\partial f^\alpha}{\partial c_K^\alpha} \quad \text{and} \quad \mu_i = \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_i^\alpha} - \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_K^\alpha}. \quad (7.13)$$

Alternatively, for the choice $\sum_{i=1}^K \mu_i = 0$, a summation over all components in the first equation of (7.9) shows that, similar to the multiplier for the sum-constraint in the phasefield case,

$$\lambda^\alpha = \frac{1}{K} \sum_{i=1}^K \frac{\partial f^\alpha}{\partial c_i^\alpha} \quad \text{and} \quad \mu_i = \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_i^\alpha} - \frac{1}{K} \sum_{i=1}^K \frac{\partial f^\alpha(\mathbf{c}, T)}{\partial c_i^\alpha}. \quad (7.14)$$

Remark 94. It is important to stress again that the optimality systems (7.9) and (7.10) a priori only apply to those phases with $h^\alpha(\boldsymbol{\phi}) \neq 0$ whereas the values for the phases with $h^\alpha(\boldsymbol{\phi}) = 0$ are a priori arbitrary as they affect neither the weighted average of the phase-specific free energies nor that of the average concentration. Nevertheless, as some of the following calculations will require the values of the \mathbf{c}^β (resp. $\tilde{\mathbf{c}}^\beta$) as well as the corresponding chemical potentials, the

¹¹Note that, in contrast to the **differences** between the derivatives, the derivatives $\frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_i^\alpha}$ themselves actually do depend upon the phase α .

standing assumption will be made that these are determined by the “natural” extension of the optimality condition above, i.e. by imposing the same optimality conditions (7.9) resp. (7.10) for all phases and not only those strictly required for actually solving the optimization problems. Despite this convention, one should still maintain the distinction between the \mathcal{P}_p and the set of all phases, since it is much more convenient to a prior only solve the actual optimality systems above and then recover the $\tilde{\mathbf{c}}^\beta$ for all phases in $\{1, \dots, N\} \setminus \mathcal{P}_p$ only when actually required from the then known value of $\boldsymbol{\mu}$ (or $\tilde{\boldsymbol{\mu}}$). \diamond

7.1.3 The Resulting Driving Force

In contrast to the simpler models based on the “common” phase-specific concentrations $\mathbf{c}^\alpha = \mathbf{c} \forall \alpha$ as in e.g. [52], the evaluation of the driving force contribution $\frac{\partial f}{\partial \phi^\alpha}$ to the phasefield equation is somewhat more complex for the model above. The primary difference is that, whereas in the former case, the phase-specific concentrations are (for a given concentration \mathbf{c}) independent upon the phasefield-values and one thus has $\frac{\partial f}{\partial \phi^\alpha} = \sum_{\beta=1}^N f^\beta(\mathbf{c}, T) \frac{\partial h^\beta}{\partial \phi^\alpha}$, the phase-specific concentration \mathbf{c}^α in the latter case depend upon the parameters of the optimization problem (7.3), i.e. $\mathbf{c}^\beta = \mathbf{c}^\beta(\phi, \mathbf{c}, T)$. This will lead to an additional contribution to $\frac{\partial f}{\partial \phi^\alpha}$ related to $\frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha}$.

More precisely, one now has¹²

$$\frac{\partial f(\phi, \mathbf{c}, T)}{\partial \phi^\alpha} = \frac{\partial}{\partial \phi^\alpha} \sum_{\beta=1}^N f^\beta(\mathbf{c}^\beta(\phi, \mathbf{c}, T), T) h^\beta(\phi) = \sum_{\beta=1}^N f^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} + \frac{\partial f^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi).$$

While this would seem to require the ability to evaluate $\frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha}$, the variational nature of the system implicitly defining the \mathbf{c}^β actually allows dispensing of this difficulty. In fact, by construction of the \mathbf{c}^β , one has $\frac{\partial f}{\partial \mathbf{c}^\beta} = \boldsymbol{\mu} + \lambda^\beta \mathbf{e}$ in the solution of Equation (7.9) and thus

$$\sum_{\beta=1}^N \frac{\partial f^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi) = \sum_{\beta=1}^N (\boldsymbol{\mu} + \lambda^\beta \mathbf{e}) \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi).$$

The first factor $\boldsymbol{\mu}$ does not depend upon the phase β though, and can thus be extracted from the sum. The second contribution $\lambda^\beta \mathbf{e} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha}$ in contrast vanishes automatically for each phase since $\mathbf{e} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} = \frac{\partial \mathbf{e} \cdot \mathbf{c}^\beta}{\partial \phi^\alpha} = \frac{\partial 1}{\partial \phi^\alpha} = 0$ by the sum-constraint on the phase-specific concentrations. Making use of the product rule and the constraint on the concentration average, one can then in a second step actually eliminate the $\frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha}$ since

$$\begin{aligned} \sum_{\beta=1}^N \frac{\partial f^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi) &= \boldsymbol{\mu} \cdot \left(\sum_{\beta=1}^N \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi) \right) = \boldsymbol{\mu} \cdot \left(\frac{\partial \left(\sum_{\beta=1}^N \mathbf{c}^\beta h^\beta(\phi) \right)}{\partial \phi^\alpha} - \sum_{\beta=1}^N \mathbf{c}^\beta \frac{\partial h^\beta(\phi)}{\partial \phi^\alpha} \right) \\ &= \boldsymbol{\mu} \cdot \left(\frac{\partial \mathbf{c}}{\partial \phi^\alpha} - \sum_{\beta=1}^N \mathbf{c}^\beta \frac{\partial h^\beta(\phi)}{\partial \phi^\alpha} \right) = -\boldsymbol{\mu} \cdot \left(\sum_{\beta=1}^N \mathbf{c}^\beta \frac{\partial h^\beta(\phi)}{\partial \phi^\alpha} \right) \end{aligned}$$

and thus one actually has, readding the arguments for clarity,

$$\frac{\partial f(\phi, \mathbf{c}, T)}{\partial \phi^\alpha} = \sum_{\beta=1}^N \left(f^\beta(\mathbf{c}^\beta(\phi, \mathbf{c}, T), T) - \boldsymbol{\mu}(\phi, \mathbf{c}, T) \cdot \mathbf{c}^\beta(\phi, \mathbf{c}, T) \right) \frac{\partial h^\beta}{\partial \phi^\alpha}. \quad (7.15)$$

A similar calculation can also, though in a slightly simpler form, be conducted for the case of the reduced formulation based on $\tilde{\mathbf{c}}$. A differentiation with respect to ϕ^α taking into account the

¹²Here the convention in Remark 94 has to be kept in mind if $\frac{\partial h^\beta}{\partial \phi^\alpha}$ does not vanish for $h^\beta(\phi) = 0$, such as for the simplest interpolation function $h^\beta(\phi) = \phi^\beta$.

dependence of the $\tilde{\mathbf{c}}^\alpha$ on ϕ first leads to $\frac{\partial \tilde{f}(\phi, \tilde{\mathbf{c}}, T)}{\partial \phi^\alpha} = \frac{\partial}{\partial \phi^\alpha} \sum_{\beta=1}^N \tilde{f}^\beta(\tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T), T) h^\beta(\phi)$. By the characterization of the $\tilde{\mathbf{c}}^\beta$ in Equation (7.10), one has on the one hand $\frac{\partial \tilde{f}^\beta}{\partial \tilde{\mathbf{c}}^\beta}(\tilde{\mathbf{c}}^\beta, T) = \tilde{\boldsymbol{\mu}}$ for all β , such that $\tilde{\boldsymbol{\mu}}$ can again be moved out of the summation. On the other hand, $\sum_{\beta=1}^N \tilde{\mathbf{c}}^\beta h^\beta(\phi) = \tilde{\mathbf{c}}$, which can, as above, be combined with the product rule to eliminate the $\frac{\partial \tilde{\mathbf{c}}^\beta}{\partial \phi^\alpha}$ making use of

$$\tilde{\boldsymbol{\mu}} \cdot \sum_{\beta=1}^N \frac{\partial \tilde{\mathbf{c}}^\beta}{\partial \phi^\alpha} h^\beta(\phi) = \tilde{\boldsymbol{\mu}} \cdot \left(\frac{\partial (\sum_{\beta=1}^N \tilde{\mathbf{c}}^\beta h^\beta(\phi))}{\partial \phi^\alpha} - \tilde{\mathbf{c}}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} \right) = \tilde{\boldsymbol{\mu}} \cdot \left(\frac{\partial \tilde{\mathbf{c}}}{\partial \phi^\alpha} - \tilde{\mathbf{c}}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} \right) = -\tilde{\boldsymbol{\mu}} \cdot \sum_{\beta=1}^N \tilde{\mathbf{c}}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$$

thus resulting in the analogue

$$\frac{\partial \tilde{f}(\phi, \tilde{\mathbf{c}}, T)}{\partial \phi^\alpha} = \sum_{\beta=1}^N \left(\tilde{f}^\beta(\tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T), T) - \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T) \right) \frac{\partial h^\beta}{\partial \phi^\alpha} \quad (7.16)$$

of Equation (7.15).

7.1.4 The Relation with the Grand Chemical Potential

Considering for simplicity the reduced formulation, it can be noted that since $\tilde{\boldsymbol{\mu}} = \frac{\partial \tilde{f}^\beta}{\partial \tilde{\mathbf{c}}^\beta}$ (see the discussion below), the contribution of each phase to the driving force in equation (7.16) corresponds precisely to the grand potential energy densities $\tilde{\Psi}^\beta$ evaluated in $\tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T)$. As already indicated in Remark 91, this is a point which has been amply discussed in various forms in the literature (see e.g. [42], [56] and [19]).

One point of view - which is in particular stressed in [56] and [19] - is that it is also possible to obtain the same driving force without having to go through the somewhat more complex dependence in terms of the phase-specific concentrations. This is achieved by replacing the free energy density with the grand potential energy density $\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) := \sum_{\alpha=1}^N \tilde{\Psi}^\alpha(\tilde{\boldsymbol{\mu}}, T) h^\alpha(\phi)$ with $\tilde{\boldsymbol{\mu}}$ considered as an independent variable (instead of, as above, as $\tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T)$), which suggests using the (total) grand potential energy density $\tilde{\Omega}_e(\phi, \tilde{\boldsymbol{\mu}}, T)$ as the relevant potential instead of the free energy density \mathcal{F} . In contrast, both [42] and [25] do not actually introduce the grand potential energy density directly, but simply state that the dependence on the phase-specific concentrations does, as above, lead to this type of driving force for minimizing the (original) free energy density.

Since both approaches result in the same driving force and the evolution of either $\tilde{\boldsymbol{\mu}}$ or $\tilde{\mathbf{c}}$ is chosen such that, up to potential differences in the choice of diffusion coefficients¹³, $\tilde{\mathbf{c}}$ satisfies a locally conservative gradient flow driven by the gradient of the (common) chemical potential, they will in fact lead to the same result. Nevertheless, the use of two a priori quite different potentials raises a number of questions which would not seem to have been discussed in sufficient detail in the literature. The following discussion will therefore try to clarify the link - and the differences in interpretation - between these two approaches.

One can first observe - this being the standard argument underlying the use of the Lagrangian in Equation (7.8) - that the optimization problem defining $\tilde{f}(\phi, \tilde{\mathbf{c}}, T)$ can equivalently be rewritten as the unconstrained minimax-problem

$$\tilde{f}(\phi, \tilde{\mathbf{c}}, T) = \min_{(\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N}} \sup_{\tilde{\boldsymbol{\mu}}} \left\{ \sum_{\alpha=1}^N \tilde{f}^\alpha(\phi, \tilde{\mathbf{c}}^\alpha, T) h^\alpha(\phi) - \tilde{\boldsymbol{\mu}} \cdot \left(\sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) - \tilde{\mathbf{c}} \right) \right\}.$$

As duality holds here by the strict convexity of the f^α , one may exchange the order of the two

¹³Their choice not being a matter of the underlying functional.

extremum operations, i.e. one can also write

$$\begin{aligned}\tilde{f}(\phi, \tilde{c}, T) &= \sup_{\tilde{\mu}} \min_{(\tilde{c}^\alpha)_{1 \leq \alpha \leq N}} \left\{ \sum_{\alpha=1}^N \tilde{f}^\alpha(\phi, \tilde{c}^\alpha, T) h^\alpha(\phi) - \tilde{\mu} \cdot \left(\sum_{\alpha=1}^N \tilde{c}^\alpha h^\alpha(\phi) - \tilde{c} \right) \right\} \\ &= \sup_{\tilde{\mu}} \left\{ \min_{(\tilde{c}^\alpha)_{1 \leq \alpha \leq N}} \left\{ \sum_{\alpha=1}^N \left(\tilde{f}^\alpha(\phi, \tilde{c}^\alpha, T) - \tilde{\mu} \cdot \tilde{c}^\alpha \right) h^\alpha(\phi) \right\} + \tilde{\mu} \cdot \tilde{c} \right\}.\end{aligned}$$

The inner minimization problem decomposes into $\sum_{\alpha=1}^N \min_{\tilde{c}^\alpha} \{ f^\alpha(\tilde{c}^\alpha, T) - \tilde{\mu} \cdot \tilde{c}^\alpha \} h^\alpha(\phi)$, where each of the separate minimization problems in fact corresponds precisely to the definition of $\Psi^\alpha(\tilde{\mu}, T)$, and thus $\tilde{f}(\phi, \tilde{c}, T) = \sup_{\tilde{\mu}} \{ \sum_{\alpha=1}^N \tilde{\Psi}^\alpha(\tilde{\mu}, T) h^\alpha(\phi) + \tilde{\mu} \cdot \tilde{c} \}$.

There are two things which can be noted here. Firstly, in the interior sum, **both** $\tilde{\mu}$ and \tilde{c} are actually independent parameters, and in particular the $\tilde{\Psi}^\alpha$ are indeed the phase-specific grand chemical potential densities $\tilde{\Psi}^\alpha(\tilde{\mu}^\alpha, T)$ evaluated for $\tilde{\mu}^\alpha = \tilde{\mu}$. Secondly, while this makes appear the average grand chemical potential density

$$\tilde{\Psi}(\phi, \tilde{\mu}, T) := \sum_{\alpha=1}^N \tilde{\Psi}^\alpha(\tilde{\mu}, T) h^\alpha(\phi), \quad (7.17)$$

it can clearly not be used to simply replace the free energy density \tilde{f} since, on the one hand, there is the additional term $\tilde{\mu} \cdot \tilde{c}$, and on the other hand, there is a remaining extremum problem in terms of $\tilde{\mu}$.

What follows though (see Remark 24 and Equation (5.19) in particular) is that the “quantitative” free energy density \tilde{f} as defined above is in fact given by the (inverse) Legendre transform of $\tilde{\Psi}$ with respect to $\tilde{\mu}$. More precisely, denoting (in analogy to Remark 24) the partial forward transform defining the $\tilde{\Psi}^\alpha$ by $\hat{\mathcal{L}}_{\tilde{c}}(f^\alpha)(\tilde{\mu}^\alpha, T)$, the free energy density can also, with the analogous notation, be characterized as

$$\tilde{f}(\phi, \tilde{c}, T) = \hat{\mathcal{L}}_{\tilde{\mu}}^{-1}[\tilde{\Psi}](\phi, \tilde{c}, T) := \sup_{\tilde{\mu}} \{ \tilde{\Psi}(\phi, \tilde{\mu}, T) + \tilde{\mu} \cdot \tilde{c} \}, \quad (7.18)$$

and is thus the inverse Legendre transform (w.r.t. $\tilde{\mu}$) of the average grand chemical potential density for $\tilde{\mu}^\alpha = \tilde{\mu}$.

Remark 95. Note that this relation - while intuitively pleasing based on the standard thermodynamic relations - is neither quite standard nor entirely obvious based upon the initial definitions. Whereas the direct definition of $\tilde{\Psi}$ involves a simple averaging procedure over the phase-specific densities $\tilde{\Psi}^\alpha(\tilde{\mu}, T)$ with a single parameter $\tilde{\mu}$, \tilde{f} is defined through the minimization of an average based on N phase-specific concentrations subject to an additional side condition on their average. In contrast, the $\tilde{\Psi}^\alpha$ themselves are a priori defined in terms of an **unconstrained** parameterized optimization in each of the \tilde{c} , $\tilde{\Psi}^\alpha(\tilde{\mu}) = \inf_{\tilde{c}^\alpha} \{ f^\alpha(\tilde{c}^\alpha) - \tilde{\mu} \cdot \tilde{c}^\alpha \}$.

That the “effective” densities \tilde{f} and $\tilde{\Psi}$ are nevertheless related in terms of a standard transform is thus not automatic but instead hinges on the internal consistency of the two definitions in terms of the common multiplier $\tilde{\mu}$. More examples of similar relations will be given in Section 7.1.6 when considering the non-isothermal case. \diamond

From this, the original optimization problem for $\tilde{\mathcal{F}}_{\tilde{c}}(\phi, \tilde{c}, T)$ is then equivalent to

$$\min_{(\phi, \tilde{c})} \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + \hat{\mathcal{L}}_{\tilde{\mu}}^{-1}[\tilde{\Psi}](\phi, \tilde{c}, T) \, d\mathbf{x} \quad (7.19)$$

and thus does clearly **not** correspond to the minimization of the grand potential energy

$$\Omega_\epsilon(\phi, \tilde{\mu}, T) := \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + \tilde{\Psi}(\phi, \tilde{\mu}, T) \, d\mathbf{x}. \quad (7.20)$$

That the driving force of the phasefield equation based on Equation (7.18) is nevertheless of the form $\sum_{\beta=1}^N \tilde{\Psi}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$ above is - similarly to appearance of $\frac{\partial f}{\partial \phi}$ despite the maximization of the entropy due the expression $\frac{\partial s}{\partial \phi}(\phi, \mathbf{c}, e) = -\frac{1}{T} \frac{\partial f}{\partial \phi}(\phi, \mathbf{c}, T)$ discussed in Section 3.2 - a result of the variational characterization of this inverse transform. More precisely, the value $\tilde{\boldsymbol{\mu}}$ for which the supremum in Equation (7.18) is actually achieved¹⁴ is characterized by

$$\frac{\partial}{\partial \tilde{\boldsymbol{\mu}}} (\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) + \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}}) = \frac{\partial \tilde{\Psi}}{\partial \tilde{\boldsymbol{\mu}}}(\phi, \tilde{\boldsymbol{\mu}}, \tilde{\mathbf{c}}) + \tilde{\mathbf{c}} \stackrel{!}{=} \mathbf{0}, \quad (7.21)$$

defining $\tilde{\boldsymbol{\mu}}$ realizing the supremum as a function of the remaining parameters, $\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T)$. Based on this maximizer, $\tilde{f}(\phi, \tilde{\mathbf{c}}, T)$ can then be rewritten explicitly in terms of $\tilde{\Psi}$ as

$$\begin{aligned} \tilde{f}(\phi, \tilde{\mathbf{c}}, T) &= \min_{\{(\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha h^\alpha(\phi) = \tilde{\mathbf{c}}\}} \left\{ \sum_{\alpha=1}^N f^\alpha(\tilde{\mathbf{c}}^\alpha, T) h^\alpha(\phi) \right\} \\ &= \Psi(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) + \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}, \end{aligned} \quad (7.22)$$

from which the driving force follows as

$$\frac{\partial \tilde{f}}{\partial \phi^\alpha}(\phi, \tilde{\mathbf{c}}, T) = \frac{\partial \Psi}{\partial \phi^\alpha}(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) + \frac{\partial}{\partial \tilde{\boldsymbol{\mu}}} (\Psi(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) + \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}) \cdot \frac{\partial \tilde{\boldsymbol{\mu}}}{\partial \phi^\alpha}.$$

By the optimality condition on $\tilde{\boldsymbol{\mu}}$, the second term drops out, leaving

$$\frac{\partial \tilde{f}}{\partial \phi^\alpha}(\phi, \tilde{\mathbf{c}}, T) = \frac{\partial \Psi}{\partial \phi^\alpha}(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) = \sum_{\beta=1}^N \tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) \frac{\partial h^\beta}{\partial \phi^\alpha},$$

and thus an expression which, together with the relation $\tilde{\Psi}^\beta = f^\beta - \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}}$ of $\tilde{\Psi}^\beta$, is the same as in Equation (7.16)¹⁵.

By the expression of the driving force in Equation (7.16) combined with the fact that the value (but not the function itself!) of $\tilde{f}^\beta(\tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T), T) - \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T)$ indeed coincides with $\tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}, T)$ evaluated for $\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T)$, it is obvious that the same result is also obtained by a straight-forward differentiation of $\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T)$ with $\tilde{\boldsymbol{\mu}}$ considered as a fixed parameter (i.e. neglecting all potential dependencies). This simple mechanism for obtaining this driving force is also the basis for the ‘‘phenomenological’’ approach taken in [56] and [19]. In fact, their approach consists in postulating for ϕ to follow a (minimizing) gradient-flow for the grand potential energy

¹⁴Note that this is a strictly concave maximization problem in $\tilde{\boldsymbol{\mu}}$.

¹⁵Before drawing this conclusion, there is one (somewhat hidden) point which remains to be verified, namely that the values of $\tilde{\boldsymbol{\mu}}$ in both equations as well as the values of $\tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}, T)$ above and the $\tilde{f}^\beta(\tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T), T) - \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}^\beta(\phi, \tilde{\mathbf{c}}, T)$ in Equation (7.16) in fact coincide as both are a priori defined based on different conditions (one as the (free) maximizer of $\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) + \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}}$, the other one as the multiplier for the (constrained) optimality system (7.10)). This now follows mostly from more standard ‘‘phase-specific’’ thermodynamics, more precisely in the (smooth convex-concave) analogue of the third point in Proposition 3 in the case of full duality, namely that $\tilde{\boldsymbol{\mu}} = \frac{\partial \tilde{f}^\beta(\tilde{\mathbf{c}}^\beta, T)}{\partial \tilde{\mathbf{c}}^\beta}$ iff $\tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}, T) = \tilde{f}^\beta(\tilde{\mathbf{c}}^\beta, T) - \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}}^\beta$ iff $\tilde{\mathbf{c}}^\beta = \frac{\partial \tilde{\Psi}^\beta}{\partial \tilde{\boldsymbol{\mu}}}$. Since this implies that simply **defining** $\tilde{\mathbf{c}}^\beta(\tilde{\boldsymbol{\mu}}, T) := \frac{\partial \tilde{\Psi}^\beta}{\partial \tilde{\boldsymbol{\mu}}}(\tilde{\boldsymbol{\mu}}, T)$ (now as a function of $\tilde{\boldsymbol{\mu}}$) with $\tilde{\boldsymbol{\mu}}$ as in Equation (7.21), one has, on the one hand, $\tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}, T) = \tilde{f}^\beta(\tilde{\mathbf{c}}^\beta(\tilde{\boldsymbol{\mu}}, T), T) - \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}}^\beta(\tilde{\boldsymbol{\mu}}, T)$ as well as $\tilde{\boldsymbol{\mu}} = \frac{\partial \tilde{f}^\beta(\tilde{\mathbf{c}}^\beta, T)}{\partial \tilde{\mathbf{c}}^\beta}(\tilde{\mathbf{c}}^\beta(\tilde{\boldsymbol{\mu}}, T), T)$ and on the other hand by the optimality condition in Equation (7.21) that

$$\tilde{\mathbf{c}} = -\frac{\partial \tilde{\Psi}}{\partial \tilde{\boldsymbol{\mu}}} = -\sum_{\alpha=1}^N \frac{\partial \tilde{\Psi}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}(\tilde{\boldsymbol{\mu}}, T) h^\alpha(\phi) = \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha(\tilde{\boldsymbol{\mu}}, T) h^\alpha(\phi),$$

i.e. $\tilde{\boldsymbol{\mu}}$ in combination with the concentration values derived from it actually satisfies the optimality system (7.10) and thus by unicity for this system that both expressions actually coincide.

functional $\Omega_\epsilon(\phi, \tilde{\boldsymbol{\mu}}, T)$ in Equation (7.20). This is combined with the standard evolution equation (here using the notation similar to [19]¹⁶ together with the Einstein summation convention) for the concentration in the isothermal case,

$$\frac{\partial c_i}{\partial t} = \nabla \cdot \left(M_{ijk} \frac{\partial \mu_j}{\partial x_k} \right), \quad i = 1, \dots, K-1 \quad \text{resp.} \quad \frac{\partial \tilde{c}}{\partial t} = \nabla \cdot (\mathbf{M} : \nabla \tilde{\boldsymbol{\mu}}) \quad (7.23)$$

for the independent concentration values.

If one takes, as proposed in both works and as the use of the grand potential energy would seem to suggest, $\tilde{\boldsymbol{\mu}}$ instead of \tilde{c} as the “fundamental” dynamic variable, this makes, together with the condition on the equilibrium of the phasespecific chemical potentials, \tilde{c} a function of $(\phi, \tilde{\boldsymbol{\mu}}, T)$, where $\tilde{c}(\phi, \tilde{\boldsymbol{\mu}}, T) = \sum_{\alpha=1}^N \tilde{c}^\alpha(\tilde{\boldsymbol{\mu}}, T) h^\alpha(\phi)$. Together with the chain-rule and the given evolution equation for ϕ , this permits rewriting the evolution equation for the reduced concentration as one for $\tilde{\boldsymbol{\mu}}$ through $\frac{\partial \tilde{c}}{\partial t} = \frac{\partial \tilde{c}}{\partial \phi} \cdot \frac{\partial \phi}{\partial t} + \frac{\partial \tilde{c}}{\partial \tilde{\boldsymbol{\mu}}} \cdot \frac{\partial \tilde{\boldsymbol{\mu}}}{\partial t}$, from which one obtains

$$\frac{\partial \tilde{\boldsymbol{\mu}}}{\partial t} = \left(\frac{\partial \tilde{c}}{\partial \tilde{\boldsymbol{\mu}}} \right)^{-1} \cdot \left(\nabla \cdot (\mathbf{M} : \nabla \tilde{\boldsymbol{\mu}}) - \frac{\partial \tilde{c}}{\partial \phi} \cdot \frac{\partial \phi}{\partial t} \right). \quad (7.24)$$

Remark 96. Note that, together with $-\tilde{c} = \frac{\partial \tilde{\Psi}}{\partial \tilde{\boldsymbol{\mu}}}$ and the evolution equation for ϕ (e.g. from Equation (6.74)) complemented by the driving force $\frac{\partial \tilde{\Psi}}{\partial \phi^\alpha} = \sum_{\beta=1}^N \tilde{\Psi}^\beta(\tilde{\boldsymbol{\mu}}, T) \frac{\partial h^\beta}{\partial \phi^\alpha}(\phi)$, this is in principle a system dependent solely on ϕ and $\tilde{\boldsymbol{\mu}}$ since the parameter T is fixed by assumption. It would thus seem that one can completely forget about the concentration and in particular the phase-specific concentrations, or at least only use them as auxiliary quantities whenever convenient.

While this is essentially true¹⁷ in the continuous case, one does well to remember that the evolution Equation (7.24) for $\tilde{\boldsymbol{\mu}}$ is designed solely for the purpose of ensuring the (conservative) evolution of \tilde{c} as in Equation (7.23) in the discrete case.

For example, when applying - as far as possible due to the term $\frac{\partial \phi}{\partial t}$ - the simplest time-discretization in terms of an explicit Euler scheme, there are several potential pitfalls. Firstly, the evaluation of the time-derivative of ϕ requires both a new and an old value $\phi^{(n+1)}$ and $\phi^{(n)}$ for evaluating the time-derivative $\frac{\partial \phi}{\partial t}$, at least when implemented as a single-step scheme¹⁸. Based solely on Equation (7.24) and the use of an explicit scheme, it is then quite natural to firstly evaluate the matrix $\frac{\partial \tilde{c}}{\partial \tilde{\boldsymbol{\mu}}}$ as $\frac{\partial \tilde{c}}{\partial \tilde{\boldsymbol{\mu}}}(\phi^{(n)}, \tilde{\boldsymbol{\mu}}^{(n)}, T)$, and secondly to explicitly evaluate the contribution due to the time-derivative of ϕ as $\frac{\partial \tilde{c}}{\partial \phi}$ as

$$\sum_{\alpha=1}^N \sum_{\beta=1}^N \tilde{c}^\beta(\phi^{(n)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) \frac{\partial h^\beta}{\partial \phi^\alpha}(\phi^{(n)}) \frac{(\phi^\alpha)^{(n+1)} - (\phi^\alpha)^{(n)}}{\Delta t}.$$

In relation with Equation (7.23) and the known data, both of these are bad choices though. In fact, instead of discretizing the “original” evolution equation for the concentration, one would like to satisfy, still within an explicit Euler approach,

$$\frac{\tilde{c}^{(n+1)} - \tilde{c}^{(n)}}{\Delta t} = \nabla \cdot (\mathbf{M} : \nabla \tilde{\boldsymbol{\mu}}^{(n)}).$$

Treating $\tilde{\boldsymbol{\mu}}$ as the primary unknown and \tilde{c} as the secondary one, this, together with the already calculated value of $\phi^{(n+1)}$ (and ∇_D denoting a discretized version of the ∇ -operator), leads to the equation

$$\tilde{c}(\phi^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n+1)}, T) \stackrel{!}{=} \tilde{c}(\phi^{(n)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) + \Delta t \nabla_D \cdot (\mathbf{M} : \nabla_D \tilde{\boldsymbol{\mu}}^{(n)}). \quad (7.25)$$

¹⁶Note that the formulation of Plapp in [56] is slightly different in that it uses a single density field, which can be related to the single independent component via the atomic volume, which is assumed to be the same for both species of atoms. For more details, see [56] or the final note in section I.B. in [19].

¹⁷Except potentially for initial conditions which are often expressed more naturally in terms of the concentration.

¹⁸An alternative of course consists in using $\phi^{(n)}$ and $\phi^{(n-1)}$, but is both more complex in terms of storage (when combined with an appropriate buffering scheme) and less accurate.

As all quantities except for $\tilde{\boldsymbol{\mu}}^{(n+1)}$ are known in this equation, a simple linearization leads to the more appropriate discrete equivalent of Equation (7.24), namely

$$\begin{aligned}\tilde{\mathbf{c}}(\boldsymbol{\phi}^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n+1)}, T) &\approx \tilde{\mathbf{c}}(\boldsymbol{\phi}^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) + \frac{\partial \tilde{\mathbf{c}}}{\partial \tilde{\boldsymbol{\mu}}}(\boldsymbol{\phi}^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n+1)}, T) \cdot (\tilde{\boldsymbol{\mu}}^{(n+1)} - \tilde{\boldsymbol{\mu}}^{(n)}) \\ &\stackrel{!}{=} \tilde{\mathbf{c}}(\boldsymbol{\phi}^{(n)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) + \Delta t \nabla_D \cdot (\mathbf{M} : \nabla_D \tilde{\boldsymbol{\mu}}^{(n)}).\end{aligned}$$

Further combining this with $\tilde{\mathbf{c}}(\boldsymbol{\phi}^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) = \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha(\tilde{\boldsymbol{\mu}}^{(n)}, T) h^\alpha(\boldsymbol{\phi}^{(n+1)})$ and the analogous expression $\tilde{\mathbf{c}}(\boldsymbol{\phi}^{(n)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) = \sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha(\tilde{\boldsymbol{\mu}}^{(n)}, T) h^\alpha(\boldsymbol{\phi}^{(n)})$ for $\mathbf{c}^{(n)}$, this then leads to the discrete update-rule

$$\begin{aligned}\tilde{\boldsymbol{\mu}}^{(n+1)} &:= \left(\frac{\partial \tilde{\mathbf{c}}}{\partial \tilde{\boldsymbol{\mu}}}(\boldsymbol{\phi}^{(n+1)}, \tilde{\boldsymbol{\mu}}^{(n)}, T) \right)^{-1} \\ &\cdot \left(\sum_{\alpha=1}^N \tilde{\mathbf{c}}^\alpha(\tilde{\boldsymbol{\mu}}^{(n)}, T) (h^\alpha(\boldsymbol{\phi}^{(n)}) - h^\alpha(\boldsymbol{\phi}^{(n+1)})) + \Delta t \nabla_D \cdot (\mathbf{M} : \nabla_D \tilde{\boldsymbol{\mu}}^{(n)}) \right),\end{aligned}\tag{7.26}$$

which differs from the discretization above in two important points.

Firstly, the linearization through the matrix $\frac{\partial \tilde{\mathbf{c}}}{\partial \tilde{\boldsymbol{\mu}}}$ is used only in terms of the (as of yet unknown) values of $\tilde{\boldsymbol{\mu}}$ but using the already known value of $\boldsymbol{\phi}^{(n+1)}$. In contrast, the most straightforward discretization does the same thing, but using a less appropriate linearization based on an “outdated” value of the phasefield-vector $\boldsymbol{\phi}^{(n)}$. Whereas this linearization will be second-order accurate in the case of Equation (7.26), it will therefore only be first-order accurate when using $\boldsymbol{\phi}^{(n)}$. Secondly, using the (also already known) differences between the weight-functions $h^\alpha(\boldsymbol{\phi})$ instead of their linearization together with $\frac{\boldsymbol{\phi}^{(n+1)} - \boldsymbol{\phi}^{(n)}}{\Delta t}$ is both cheaper and more accurate (actually exact) than what would result from an a priori natural discretization based on Equation (7.24).

With this choice of parameters, the resulting update rule for $\tilde{\boldsymbol{\mu}}$ clearly corresponds to a single Newton-step for the discrete update rule (7.25) for the concentration. An obvious advantage of using such a “blind” single-step procedure in the update rule (7.26) as compared to actually solving Equation (7.25) is that this is certainly cheaper and requires fewer parameters as there is for example no control of the residual before accepting the new $\tilde{\boldsymbol{\mu}}$ -value as sufficiently accurate. This is at the same time also potentially a severe disadvantage, since this essentially restricts the use of Equation (7.26) to situations where the changes in $\tilde{\boldsymbol{\mu}}$ between successive time-steps are small enough for this to be the case. A particular but practically quite useful setting where this never causes any issues is of course when using parabolic approximations of the free energy densities, as, with the $\tilde{\mathbf{c}}^\alpha$ being affine functions of $\tilde{\boldsymbol{\mu}}$, the linearization is actually exact. \diamond

Even though taking $\tilde{\boldsymbol{\mu}}$ instead of $\tilde{\mathbf{c}}$ as the primary unknown (and thus a quantity which is seemingly independent of $\boldsymbol{\phi}$) would at first sight seem to justify the use of the partial differentiation of $\tilde{\Psi}$ with respect to $\boldsymbol{\phi}$ only and thus the use of the grand potential energy density as the relevant functional, this is somewhat misleading. More precisely, from an optimization point of view, the relevance of the functional is not really a question of whether one parameterizes an intermediary evolution in terms of $\tilde{\mathbf{c}}$ or $\tilde{\boldsymbol{\mu}}$ (there being a one-to-one correspondence between the two), but of whether the final steady-state can indeed be related to a critical point of this functional or not. This is closely related to the discussion in [41]. In fact, once evolved to a steady-state, the (reduced) chemical potential $\tilde{\boldsymbol{\mu}}$ will, regardless of the choice of parameterization in $\tilde{\mathbf{c}}$ or $\tilde{\boldsymbol{\mu}}$, satisfy the equation

$$-\nabla \cdot (\mathbf{M} : \nabla \tilde{\boldsymbol{\mu}}) = \mathbf{0}$$

together with the natural isolating boundary conditions $\frac{\partial \tilde{\boldsymbol{\mu}}}{\partial \mathbf{n}} = \mathbf{0}$. This equation will (under the natural coercivity condition on \mathbf{M}) enforce for $\tilde{\boldsymbol{\mu}}$ to be equal to a constant, i.e. $\tilde{\boldsymbol{\mu}}(\mathbf{x}) = \tilde{\boldsymbol{\eta}}$ for some $(K-1)$ -dimensional vector $\tilde{\boldsymbol{\eta}}$, which is just the well-known condition of the equilibration of the

chemical potentials. As the conservative gradient flow for \tilde{c} in Equation (7.23) is, in combination with the isolating boundary conditions, specifically chosen such that the total concentration $\int_{\Omega} \tilde{c} \, d\mathbf{x}$ remains equal to its initial value \tilde{C}_0 , this final value $\tilde{\eta}$ is not arbitrary, but fixed through the condition

$$\int_{\Omega} \tilde{c}(\phi, \tilde{\eta}, T) \, d\mathbf{x} = \int_{\Omega} -\frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}}(\phi, \tilde{\eta}, T) \, d\mathbf{x} = \tilde{C}_0. \quad (7.27)$$

From this, it is obvious that, even if one initially chooses $\tilde{\mu}$ as an “independent” unknown, the final values of $\tilde{\mu}$ are of the form $\tilde{\mu} = \tilde{\eta}(\phi, \tilde{C}_0, T)$, where $\tilde{\eta}$ depends implicitly (and in a non-local fashion) on ϕ through the consistency condition (7.27) on the total concentration¹⁹.

As the (physically obviously important) dynamics in themselves ultimately have no impact on the final values of the functional, one might just as well directly consider Ω_{ϵ} using this final state $\tilde{\eta}(\phi, T)$, reducing Ω_{ϵ} to a function of ϕ and the fixed temperature alone,

$$\Omega_{\epsilon}(\phi, \tilde{\eta}(\phi, \tilde{C}_0, T), T) = \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + \tilde{\Psi}(\phi, \tilde{\eta}(\phi, \tilde{C}_0, T), T) \, d\mathbf{x}$$

together with Equation (7.27) as a side-condition. An actual minimizer for Ω_{ϵ} in terms of ϕ (being the only free parameter now), would thus have to satisfy the analogue of the steady-state phasefield equation 6.64, but with the driving force contribution due to $\tilde{\Psi}$ now derived from

$$\int_{\Omega} \frac{\partial \tilde{\Psi}}{\partial \phi} \cdot \delta \phi + \frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}} \cdot \delta \tilde{\eta}(\delta \phi, T) \, d\mathbf{x}$$

where $\delta \tilde{\eta}(\delta \phi)$ can be determined through the differentiation

$$\begin{aligned} & \int_{\Omega} \frac{\partial \tilde{c}}{\partial \phi}(\phi, \tilde{\eta}(\phi, T), T) \cdot \delta \phi + \frac{\partial \tilde{c}}{\partial \tilde{\mu}}(\phi, \tilde{\eta}(\phi, T), T) \cdot \delta \tilde{\eta}(\delta \phi) \, d\mathbf{x} \\ &= - \int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \phi \partial \tilde{\mu}}(\phi, \tilde{\eta}(\phi, T), T) \cdot \delta \phi + \frac{\partial^2 \tilde{\Psi}}{\partial \tilde{\mu}^2}(\phi, \tilde{\eta}(\phi, T), T) \cdot \delta \tilde{\eta}(\delta \phi) \, d\mathbf{x} = \mathbf{0} \end{aligned}$$

of the constraint in Equation (7.27). Since $\delta \tilde{\eta}$ is a constant vector, this equation is easily resolved explicitly for $\delta \tilde{\eta}(\delta \phi)$ as

$$\delta \tilde{\eta}(\delta \phi) = - \left(\int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \tilde{\mu}^2}(\phi, \tilde{\eta}(\phi, T), T) \, d\mathbf{x} \right)^{-1} \cdot \int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \phi \partial \tilde{\mu}}(\phi, \tilde{\eta}(\phi, T), T) \cdot \delta \phi \, d\mathbf{x},$$

showing that an actual minimizer of Ω_{ϵ} would have to be based on the driving force corresponding to

$$\int_{\Omega} \frac{\partial \tilde{\Psi}}{\partial \phi} \cdot \delta \phi - \frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}} \cdot \left(\int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \tilde{\mu}^2} \, d\mathbf{y} \right)^{-1} \cdot \int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \phi \partial \tilde{\mu}} \cdot \delta \phi \, d\mathbf{y} \, d\mathbf{x}.$$

Exchanging the order of integration in the last part, this contribution simplifies to the “local” density

$$\left(\frac{\partial \tilde{\Psi}}{\partial \phi} - \left(\int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \tilde{\mu}^2} \, d\mathbf{y} \right) \cdot \left(\int_{\Omega} \frac{\partial^2 \tilde{\Psi}}{\partial \tilde{\mu}^2} \, d\mathbf{y} \right)^{-1} \cdot \frac{\partial^2 \tilde{\Psi}}{\partial \phi \partial \tilde{\mu}} \right) \cdot \delta \phi. \quad (7.28)$$

This is obviously quite different from the expression resulting from the first part alone, showing that steady-state solutions of the system above do **not** correspond to actual minimizers (or more precisely critical points in general) of Ω_{ϵ} under the side-condition on the total concentration.

¹⁹This is essentially the same argument as in [41].

In contrast, repeating an analogous argument for the free energy functional, with $\tilde{\mathbf{c}}$ parameterized in terms of $(\phi, \tilde{\boldsymbol{\eta}}(\phi, T), T)$, leads to the condition

$$\int_{\Omega} \frac{\partial \tilde{f}}{\partial \phi}(\phi, \tilde{\mathbf{c}}(\phi, \tilde{\boldsymbol{\eta}}(\phi, \tilde{\mathbf{C}}_0, T), T), T) \cdot \delta \phi + \frac{\partial \tilde{f}}{\partial \tilde{\mathbf{c}}}(\phi, \tilde{\mathbf{c}}(\phi, \tilde{\boldsymbol{\eta}}(\phi, \tilde{\mathbf{C}}_0, T), T), T) \cdot \delta \tilde{\mathbf{c}}(\delta \phi) \, d\mathbf{x}, \quad (7.29)$$

where $\delta \tilde{\mathbf{c}}(\delta \phi)$ (despite its a priori more complex dependence on both $\delta \phi$ and $\delta \tilde{\boldsymbol{\mu}}(\delta \phi)$) can again be characterized from a differentiation of Equation (7.27) through the very simple relation $\int_{\Omega} \delta \tilde{\mathbf{c}}(\delta \phi) \, d\mathbf{x} = \mathbf{0}$. By the equilibrium condition in terms of the concentration through the equilibration of the chemical potentials, one in addition has

$$\frac{\partial \tilde{f}}{\partial \tilde{\mathbf{c}}}(\phi, \tilde{\mathbf{c}}(\phi, \tilde{\boldsymbol{\eta}}(\phi, \tilde{\mathbf{C}}_0, T), T), T) = \tilde{\boldsymbol{\eta}}(\phi, \tilde{\mathbf{C}}_0, T)$$

with $\tilde{\boldsymbol{\eta}}$ independent of \mathbf{x} . It follows that the second contribution in Equation (7.29) simply drops out and the driving force in Equation (7.16) is thus indeed consistent with a variational principle in terms of the free energy together with the side-condition given by the conservation of the total concentration $\tilde{\mathbf{C}}_0$.

This can alternatively also be seen by, instead of relying on the conservative gradient flow for $\tilde{\mathbf{c}}$ with respect to \mathcal{F}_{ϵ} , including the constraint $\int_{\Omega} \tilde{\mathbf{c}} \, d\mathbf{x} = \tilde{\mathbf{C}}_0$ directly into the problem by augmenting \mathcal{F}_{ϵ} with this constraint, i.e. by considering

$$\min_{(\phi, \tilde{\mathbf{c}})} \sup_{\tilde{\boldsymbol{\eta}}} \left\{ \mathcal{F}_{\epsilon} - \tilde{\boldsymbol{\eta}} \cdot \left(\int_{\Omega} \tilde{\mathbf{c}} \, d\mathbf{x} - \tilde{\mathbf{C}}_0 \right) \right\}.$$

Using Equation (7.18), this may also be written as

$$\min_{(\phi, \tilde{\mathbf{c}})} \sup_{\tilde{\boldsymbol{\eta}}} \left\{ \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + \sup_{\tilde{\boldsymbol{\mu}}} \left\{ \tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) + \tilde{\boldsymbol{\mu}} \cdot \tilde{\mathbf{c}} \right\} - \tilde{\boldsymbol{\eta}} \cdot \left(\int_{\Omega} \tilde{\mathbf{c}} \, d\mathbf{x} - \tilde{\mathbf{C}}_0 \right) \right\}.$$

Assuming one can interchange the minimization with respect to $\tilde{\mathbf{c}}$ with the suprema in $\tilde{\boldsymbol{\mu}}$ and $\tilde{\boldsymbol{\eta}}$, the minimization with respect to $\tilde{\mathbf{c}}$, with the only term depending on $\tilde{\mathbf{c}}$ given by $\int_{\Omega} (\tilde{\boldsymbol{\mu}} - \tilde{\boldsymbol{\eta}}) \cdot \tilde{\mathbf{c}} \, d\mathbf{x}$ leads to

$$\sup_{(\tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\mu}})} \inf_{\tilde{\mathbf{c}}} \left\{ \int_{\Omega} \tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) + (\tilde{\boldsymbol{\mu}} - \tilde{\boldsymbol{\eta}}) \cdot \tilde{\mathbf{c}} \, d\mathbf{x} + \tilde{\boldsymbol{\eta}} \cdot \tilde{\mathbf{C}}_0 \right\} = \sup_{(\tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\mu}})} \left\{ \begin{array}{ll} \int_{\Omega} \tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T) \, d\mathbf{x} + \tilde{\boldsymbol{\eta}} \cdot \tilde{\mathbf{C}}_0 & \tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\eta}} \text{ a.e.}, \\ -\infty & \text{else} \end{array} \right\}.$$

The second case will obviously be eliminated by the sup-operations, such that it suffices to consider the case with $\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\eta}}$ which can then be reduced to a global maximization operation in the vector $\tilde{\boldsymbol{\mu}}$, $\sup_{\tilde{\boldsymbol{\eta}}} \left\{ \int_{\Omega} \tilde{\Psi}(\phi, \tilde{\boldsymbol{\eta}}, T) \, d\mathbf{x} + \tilde{\boldsymbol{\eta}} \cdot \tilde{\mathbf{C}}_0 \right\}$. Combining this with the contributions by a and w , an alternative formulation of the global problem is thus given by

$$\min_{\phi} \sup_{\tilde{\boldsymbol{\eta}}} \left\{ \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + \tilde{\Psi}(\phi, \tilde{\boldsymbol{\eta}}, T) \, d\mathbf{x} + \tilde{\boldsymbol{\eta}} \cdot \tilde{\mathbf{C}}_0 \right\} = \min_{\phi} \sup_{\tilde{\boldsymbol{\eta}}} \left\{ \Omega_{\epsilon}(\phi, \tilde{\boldsymbol{\eta}}, T) + \tilde{\boldsymbol{\eta}} \cdot \tilde{\mathbf{C}}_0 \right\}. \quad (7.30)$$

Note that since the phasefield terms play no role with respect to the supremum in $\tilde{\boldsymbol{\mu}}$, this is basically just a ‘‘global’’ version of the inverse transform defining $\tilde{f}(\phi, \tilde{\mathbf{c}}, T)$ in terms of $\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}, T)$ as in Equation (7.18). More precisely, the optimality condition for $\tilde{\boldsymbol{\mu}}$ in Equation (7.21) defines $\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T)$ as a function of the given parameters $(\phi, \tilde{\mathbf{c}}, T)$. This value of $\tilde{\boldsymbol{\mu}}$ is such that, on the one hand, $\tilde{\mathbf{c}} = -\frac{\partial \tilde{\Psi}}{\partial \tilde{\boldsymbol{\mu}}}(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T)$ and on the other hand, from Equation (7.22) and the definition (7.5) of $\tilde{f}(\phi, \tilde{\mathbf{c}}, T)$, such that the value of $\tilde{\Psi}(\phi, \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T), T) + \tilde{\boldsymbol{\mu}}(\phi, \tilde{\mathbf{c}}, T) \cdot \tilde{\mathbf{c}}$ is precisely

the one of the effective free energy obtained by an energetically optimal distribution of the given average concentration onto the phase-specific ones.

In an analogous manner, the optimality condition with respect to $\tilde{\eta}$ in Equation (7.30) is given by $\tilde{\mathbf{C}}_0 = -\int_{\Omega} \frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}}(\phi, \tilde{\eta}, T) d\mathbf{x}$, i.e. an integral version of $\tilde{c} = -\frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}}$, which serves to define $\tilde{\eta} = \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T)$ with the total concentration $\tilde{\mathbf{C}}_0$ replacing the local one in $\tilde{\mu}(\phi, \tilde{c}, T)$. In addition, this value of $\tilde{\eta}$ satisfies an integral analogue of Equation (7.22), namely one has

$$\begin{aligned} \min_{\{\tilde{c}: \int_{\Omega} \tilde{c} d\mathbf{x} = \tilde{\mathbf{C}}_0\}} \left\{ \int_{\Omega} \tilde{f}(\phi, \tilde{c}, T) d\mathbf{x} \right\} &= \sup_{\tilde{\eta}} \left\{ \int_{\Omega} \tilde{\Psi}(\phi, \tilde{\eta}, T) d\mathbf{x} + \tilde{\eta} \cdot \tilde{\mathbf{C}}_0 \right\} \\ &= \int_{\Omega} \tilde{\Psi}(\phi, \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T), T) d\mathbf{x} + \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T) \cdot \tilde{\mathbf{C}}_0, \end{aligned} \quad (7.31)$$

and thus the value of $\int_{\Omega} \tilde{\Psi}(\phi, \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T), T) d\mathbf{x} + \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T) \cdot \tilde{\mathbf{C}}_0$ is the one obtained from an energetically optimal redistribution of the given concentration $\tilde{\mathbf{C}}_0$ over the domain.

The primary interest of this observation in terms of the current discussion is that it can serve to derive an “extension” of the expression for the driving force in Equation (7.16). Whereas this expression was obtained based on the local minimization of \tilde{f} in terms of the \mathbf{c}^α given the local concentration \tilde{c} , Equation (7.31) defines the global optimal contribution due to the bulk free energy density through an a priori quite complex double-minimization as

$$\min_{\{\tilde{c}: \int_{\Omega} \tilde{c} d\mathbf{x} = \tilde{\mathbf{C}}_0\}} \left\{ \left(\tilde{c}^\alpha \right)_{1 \leq \alpha \leq N} : \sum_{\alpha=1}^N \tilde{c}^\alpha h^\alpha(\phi) = \tilde{c} \right\} \left\{ \int_{\Omega} \sum_{\alpha=1}^N f^\alpha(\tilde{c}^\alpha, T) h^\alpha(\phi) d\mathbf{x} \right\}$$

enforcing an optimal redistribution of the given total concentration $\tilde{\mathbf{C}}_0$ over both the domain and the individual phases. Regardless, in the same manner as before, it follows from Equation (7.31) that the driving force contribution to this problem has the same simple structure as before. More precisely, instead of differentiating the expression for $\int_{\Omega} \tilde{f}(\phi, \tilde{c}(\phi, T), T)$ where $\tilde{c} = \tilde{c}(\phi, T)$ is the minimizer for the left-most expression in Equation (7.31), one can also differentiate the right-most one, leading to

$$\delta \left(\int_{\Omega} \tilde{f}(\phi, \tilde{c}(\phi, T), T) d\mathbf{x} \right) (\delta\phi) = \int_{\Omega} \frac{\partial \tilde{\Psi}}{\partial \phi} \cdot \delta\phi + \frac{\partial \tilde{\Psi}}{\partial \tilde{\eta}} \cdot \delta\tilde{\eta}(\delta\phi) d\mathbf{x} + \delta\tilde{\eta}(\delta\phi) \cdot \tilde{\mathbf{C}}_0$$

with $\delta\tilde{\eta}(\delta\phi)$ denoting the increment of $\tilde{\eta}(\delta\phi)$ as a (non-local) function of $\delta\phi$. As $\tilde{\eta}$ is such that $\tilde{\mathbf{C}}_0 = -\int_{\Omega} \frac{\partial \tilde{\Psi}}{\partial \tilde{\mu}}(\phi, \tilde{\eta}, T) d\mathbf{x}$ though, the last two terms drop out regardless of the value of $\delta\tilde{\eta}(\delta\phi)$, thus again reducing the derivative with respect to ϕ (despite the additional global dependence of the optimal \tilde{c} on ϕ) to the simple local expression

$$\frac{\partial \tilde{f}}{\partial \phi^\alpha}(\phi, \tilde{c}(\phi, \tilde{\mathbf{C}}_0, T), T) = \sum_{\beta=1}^N \tilde{\Psi}^\beta(\phi, \tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T), T) \frac{\partial h^\beta}{\partial \phi^\alpha}(\phi).$$

Summarizing the discussion above, it can be seen that while the $\tilde{\Psi}^\alpha$ appear very naturally as the relevant phase-specific quantities for the driving force in the phasefield equation, the link between this model and the grand potential energy functional Ω_ϵ is a somewhat tenuous one. While the use of this potential as suggested in [56] and [19] is clearly a very efficient approach for “deriving” this driving force, this should rather be considered as a purely formal device²⁰.

²⁰In particular, as seen in Equation (7.28), a more careful derivation of the driving force corresponding to the minimization of Ω_ϵ under the implicit constraint of maintaining the total concentration $\int_{\Omega} \tilde{c} d\mathbf{x}$ enforced through the evolution equation for \tilde{c} (whether expressed in $\tilde{\mu}$ or not) leads to a quite different expression. In contrast, as seen based on Equation (7.31), it is only in combination with the additional contribution by $\tilde{\eta}(\phi, \tilde{\mathbf{C}}_0, T) \cdot \tilde{\mathbf{C}}_0$ to Ω_ϵ - and thus reverting back to the free energy - that one recovers consistency of the driving force with a global minimization property.

As such this approach is in a certain sense in competition with the more common approach of “guessing” the correct equations based upon a straightforward partial differentiation of the Lagrangian

$$L(\phi, \tilde{c}, T, (\tilde{c}^\alpha)_{1 \leq \alpha \leq N}, \tilde{\mu}) = \sum_{\alpha=1}^N \tilde{f}^\alpha(\tilde{c}^\alpha, T) h^\alpha(\phi) - \tilde{\mu} \cdot \left(\sum_{\alpha=1}^N \tilde{c}^\alpha h^\alpha(\phi) - \tilde{c} \right).$$

While this approach is a priori also purely formal, obtaining the relevant driving force is then just as simple as when starting from Ω_ϵ as it suffices to partially differentiate L with respect to ϕ^α , with the additional advantage of being, even at a purely formal level, easier to link to the original problem one is trying to solve²¹.

Remark 97. It is important to note that the discussion above is only valid in this form when the concentration is (through the imposed conservative flow of \tilde{c}) subject to isolating boundary conditions and therefore implicitly the constraint on the total concentration $\int_\Omega \tilde{c} d\mathbf{x}$. If, in contrast, one chooses to enforce a constant value of the chemical potential $\hat{\mu}$ on the outer boundary, it is clear that the phasefield equation together with the driving force given in Equation (7.16), will indeed minimize the grand potential energy $\Omega_\epsilon(\phi, \hat{\mu}, T)$. In fact, since the chemical potential is, up to some potential delay through the diffusion equation for \tilde{c} , then fixed to the value $\hat{\mu}$, there will be no additional contribution as in Equation (7.28).

At the same time, there will not anymore be a natural minimization interpretation in terms of \mathcal{F} , a property which is closely related to the discussion in the earlier paper [41] by Kim, Kim and Suzuki. More precisely, the analysis in [41] investigates the influence on the minimization of the free energy $\mathcal{F}(\phi, c)$ if one, instead of treating both c and ϕ (here in a reduced scalar formulation) as independent variables, considers the steady-state concentration profile for c as a function of the steady-state phasefield profile $\phi(\mathbf{x})$. Their argument is that any equilibrium profile of the concentration has to satisfy the equilibration of the chemical potentials μ with the prescribed one, i.e. $\frac{\partial f}{\partial c}(c, \phi) = \hat{\mu} = \text{const}$ and thus implicitly defines c as a function of ϕ (and $\hat{\mu}$). They then continue to argue that defining $F(c(\phi), \phi) := f(c(\phi), \phi) - \hat{\mu}c(\phi)$ and thus such that F satisfies $\frac{\partial F}{\partial c} = \frac{\partial f}{\partial c} - \hat{\mu} = 0$ and $\frac{dF}{d\phi} = \frac{\partial F}{\partial \phi} + \frac{\partial F}{\partial c} \frac{\partial c}{\partial \phi} = \frac{\partial F}{\partial \phi} = \frac{\partial f}{\partial \phi}(c(\phi), \phi)$, the driving force $\frac{\partial f}{\partial \phi}$ for the phasefield equation can equivalently be reformulated in terms of $\frac{dF}{d\phi}$.

While this argument is correct, it needs to be interpreted with some care in the present context²². In particular, the analysis in [41] should **not** be considered in the sense that the driving force for a minimization of the free energy in this case should be given by the derivative of $f(c, \phi) - \hat{\mu}c$. It is thus quite different from the analysis in the paper [42] by the same authors, where, despite the fact that it is still based on the minimization of the free energy, the use of two phase-specific concentration fields defined as functions of (ϕ, c) in terms of the equality of the phase-specific (reduced) chemical potentials does in fact lead to the driving force being given by the differences of $f^\alpha(c^\alpha) - \mu c^\alpha$ ²³.

Instead, it simply shows a variational inconsistency in the sense that under the constraint in terms of $\hat{\mu}$, using the driving force $\frac{\partial f}{\partial \phi}$ derived as a partial derivative from the free energy density, f will in fact not minimize the free energy but the grand chemical potential (i.e. the

²¹I.e. in particular, when the \tilde{c}^α satisfy the sum-constraint, L reduces to $\sum_{\alpha=1}^N f^\alpha(\tilde{c}^\alpha, T) h^\alpha(\phi)$ and equating the partial differentiations with respect to $\tilde{\mu}$ and \tilde{c}^α to $\mathbf{0}$ leads to the constraint and optimality condition $\frac{\partial f^\alpha}{\partial \tilde{c}^\alpha} = \tilde{\mu}$ respectively.

²²In [41], the appearance of this total differential is primarily used to enable a more standard analysis of the energetics of the transition region in the one-dimensional case.

²³In contrast to the former paper, this is based on a purely local analysis and does not consider any additional dependence of the average concentration c on ϕ .

quantity whose total derivative $\frac{\partial f}{\partial \phi}$ represents in combination with the constraint). As such, it is essentially a restatement of the well-known fact that the minimization of the free energy is not the appropriate variational principle when dealing with an open system.

In contrast, actually minimizing the free energy under the constraint on the chemical potential would require using the total derivative of f , i.e.

$$\frac{df}{d\phi}(\phi, c(\phi)) = \frac{\partial f}{\partial \phi} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial \phi} = \frac{\partial f}{\partial \phi} + \hat{\mu} \frac{\partial c}{\partial \phi},$$

which could then, using the same argument as in [41], be rewritten as $\frac{\partial f}{\partial \phi} - \hat{\mu} \left(\frac{\partial^2 f}{\partial c^2} \right)^{-1} \frac{\partial^2 f}{\partial \phi \partial c}$. \diamond

Remark 98. It is quite interesting to observe what happens if the analysis in [42] is combined with the one in [41] by again enforcing a **prescribed** chemical potential²⁴. Firstly, adjusting the previous argument to the notation above, it was already observed that $\frac{\partial \tilde{f}}{\partial \phi} = \sum_{\beta=1}^N \tilde{\Psi}^\beta(\hat{\mu}, T) \frac{\partial h^\beta}{\partial \phi^\alpha}$.

Since this is, for a fixed value of $\hat{\mu}$, clearly the same as $\frac{d\tilde{\Psi}}{d\phi}$, using $\frac{\partial f}{\partial \phi}$ as the driving force does in fact lead to the minimization of Ω_ϵ . Secondly, the total derivative of $\tilde{f}(\phi, \tilde{c}(\phi, \hat{\mu}, T), T)$ is then given by

$$\frac{d\tilde{f}}{d\phi} = \frac{\partial \tilde{f}}{\partial \phi} + \frac{\partial \tilde{f}}{\partial \tilde{c}} \cdot \frac{\partial \tilde{c}}{\partial \phi}.$$

As in this model the average concentration \tilde{c} is by construction given by the weighted average $\tilde{c} = \sum_{\alpha=1}^N \tilde{c}(\hat{\mu}, T) h^\alpha(\phi)$ of the \tilde{c}^α corresponding to the - then given - chemical potential $\hat{\mu}$, there is in fact no need to invoke an implicit function theorem for determining $\frac{\partial \tilde{c}}{\partial \phi}$ as these derivatives can directly be obtained from this definition. Combining this with the expression for $\frac{\partial \tilde{f}}{\partial \phi^\alpha}$ and the chain-rule for the differentiation of \tilde{f} , one thus obtains

$$\frac{d\tilde{f}}{d\phi^\alpha} = \sum_{\beta=1}^N \tilde{\Psi}^\beta(\hat{\mu}, T) \frac{\partial h^\beta}{\partial \phi^\alpha} + \hat{\mu} \cdot \sum_{\beta=1}^N \tilde{c}^\beta(\hat{\mu}, T) \frac{\partial h^\beta}{\partial \phi^\alpha} = \sum_{\beta=1}^N f^\beta(\tilde{c}^\beta(\hat{\mu}, T), T) \frac{\partial h^\beta}{\partial \phi^\alpha}.$$

In contrast to the minimization of \mathcal{F} subject to a fixed total concentration, the correct driving force for the minimization of \mathcal{F} subject to a fixed chemical potential is therefore indeed given in terms of the \tilde{f}^α . \diamond

7.1.5 The Practical Evaluation of $f(\phi, c, T)$ and the Chemical Potential

Given that f is now only implicitly defined in terms of either the abstract minimization problem Equation (7.3) (resp. through its alternative representation \tilde{f} in Equation (7.5)) or the more concrete optimality system in Equation (7.9) (resp. in Equation (7.10)), it is clear that the evaluation of the free energy function is more involved than for an explicit formulation in terms of the phase-averaged concentrations c as in e.g. [52]. Therefore, each evaluation of $f(\phi, c, T)$ now requires solving for the $(c^\alpha)_{\alpha \in \mathcal{P}_p}$ realizing the minimum in Equation (7.3). As the $f^\alpha(c^\alpha, T)$ are in general nonlinear (though strictly convex) functions of the phase-specific concentrations c^α , this has to be done - except for particularly simple examples such as an ideal solution model - using some iterative solution approach. The natural choice here is of course a Newton-type scheme, based on either Equation (7.9) or Equation (7.10).

The Reduced Formulation

Beginning with the somewhat simpler case of the reduced formulation, each Newton step consists in, given the current estimates of $(\tilde{c}^\alpha)_{\alpha \in \mathcal{P}_p}^{(n)}$ (and therefore trivially also $(c_K^\alpha)^{(n)} = 1 -$

²⁴Given that both reduced chemical potentials are enforced to be equal, the only consistent choice for them is to be equal to the prescribed one.

$\sum_{i=1}^{K-1} (c_i^\alpha)^{(n)}$, determining the updated concentration values $(\tilde{c}^\alpha)^{(n+1)} = ((\tilde{c}^\alpha)^{(n)} + \delta\tilde{c}^\alpha)_{\alpha \in \mathcal{P}_p}$ and the associated updated multiplier $\tilde{\mu}^{(n+1)}$ by either setting $\tilde{\mu}^{(n)} + \delta\tilde{\mu}$ and solving the linearized problem

$$\begin{cases} \tilde{\mu}^\alpha((\tilde{c}^\alpha)^{(n)}, T) + \frac{\partial \tilde{\mu}^\alpha}{\partial \tilde{c}^\alpha}((\tilde{c}^\alpha)^{(n)}, T) \delta\tilde{c}^\alpha - (\tilde{\mu}^{(n)} + \delta\tilde{\mu}) = \mathbf{0}, \\ \sum_{\alpha \in \mathcal{P}_p} ((\tilde{c}^\alpha)^{(n)} + \delta c^\alpha) h^\alpha(\phi) = \tilde{c}. \end{cases} \quad (7.32)$$

Alternatively, as the original equations $\tilde{\mu}^\alpha \stackrel{!}{=} \tilde{\mu}$ and thus automatically also (7.32) are in fact linear in $\tilde{\mu}$, one can also directly use $\tilde{\mu}^{(n+1)}$ and solve

$$\begin{cases} \tilde{\mu}^\alpha((\tilde{c}^\alpha)^{(n)}, T) + \frac{\partial \tilde{\mu}^\alpha}{\partial \tilde{c}^\alpha}((\tilde{c}^\alpha)^{(n)}, T) \delta\tilde{c}^\alpha - \tilde{\mu}^{(n+1)} = \mathbf{0}, \\ \sum_{\alpha \in \mathcal{P}_p} ((\tilde{c}^\alpha)^{(n)} + \delta c^\alpha) h^\alpha(\phi) = \tilde{c}. \end{cases} \quad (7.33)$$

Remark 99. It is clear that²⁵ both the incremental (in $\tilde{\mu}$) System (7.32) and the non-incremental System (7.33) will lead to the same value for $\tilde{\mu}^{(n+1)}$. Nevertheless, using an incremental formulation can simplify certain modifications (such as the use of a damped Newton scheme when the convexity is lost due to additional contributions) and, as pointed out below, allows a simpler reuse of the algorithm when determining the c^α for a **given** $\tilde{\mu}$. It will therefore be preferred in the following. If desired, the non-incremental version can easily be obtained from the description below by simply replacing $\delta\mu$ with $\tilde{\mu}^{(n+1)}$ and the old value $\tilde{\mu}^{(n)}$ by $\mathbf{0}$. \diamond

Assuming for notational simplicity that $\mathcal{P}_p = \{1, 2, \dots, N_p\}$, the matrix-vector form of the system (7.32) is given by²⁶

$$\begin{pmatrix} \left(\begin{array}{cccc} \left(\frac{\partial^2 \tilde{f}^1}{\partial (\tilde{c}^1)^2} \right) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \left(\frac{\partial^2 \tilde{f}^2}{\partial (\tilde{c}^2)^2} \right) & \dots & \mathbf{0} \\ & \ddots & \ddots & \ddots \\ \mathbf{0} & \dots & \mathbf{0} & \left(\frac{\partial^2 \tilde{f}^{N_p}}{\partial (\tilde{c}^{N_p})^2} \right) \end{array} \right) & \begin{pmatrix} -\mathbf{I} \\ -\mathbf{I} \\ \vdots \\ -\mathbf{I} \end{pmatrix} & \begin{pmatrix} \delta\tilde{c}^1 \\ \delta\tilde{c}^1 \\ \vdots \\ \delta\tilde{c}^{N_p} \\ \delta\tilde{\mu} \end{pmatrix} & = & \begin{pmatrix} \mathbf{r}_{\tilde{\mu}^1}^{(n)} \\ \mathbf{r}_{\tilde{\mu}^2}^{(n)} \\ \vdots \\ \mathbf{r}_{\tilde{\mu}^{N_p}}^{(n)} \\ \mathbf{r}_{\tilde{c}}^{(n)} \end{pmatrix} \end{pmatrix},$$

where the residuals are given by

$$\mathbf{r}_{\tilde{\mu}^\alpha}^{(n)} := \tilde{\mu}^{(n)} - (\tilde{\mu}^\alpha)^{(n)} \quad \text{and} \quad \mathbf{r}_{\tilde{c}}^{(n)} := \mathbf{c} - \sum_{\alpha=1}^{N_p} (\tilde{c}^\alpha)^{(n)} h^\alpha(\phi) \quad (7.34)$$

and the dependence of the $\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2}$ on $((\tilde{c}^\alpha)^{(n)}, T)$ has been suppressed. A noteworthy feature of this system is that the upper left part is block-diagonal as each \tilde{f}^α only depends on the phase-specific concentrations of this particular phase, and each $(K-1) \times (K-1)$ -block is formed by an s.p.d. matrix due to the convexity of the \tilde{f}^α w.r.t. \tilde{c}^α . Due to this particular structure, the system can conveniently be solved using a Schur complement approach, i.e. by first eliminating the $\delta\tilde{c}^\alpha$ in terms of the $\mathbf{r}_{\tilde{\mu}^\alpha}^{(n)}$ and $\delta\tilde{\mu}$ (resp. $\tilde{\mu}^{(n+1)}$) as

$$\delta\tilde{c}^\alpha = \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} (\mathbf{r}_{\tilde{\mu}^\alpha}^{(n)} + \delta\tilde{\mu}) = \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} (\tilde{\mu}^{(n+1)} - (\tilde{\mu}^\alpha)^{(n)}) \quad \alpha \in \mathcal{P}_p. \quad (7.35)$$

²⁵ Also see Remark 101 below.

²⁶ Note that this system in this form does not correspond to a classical saddle point matrix. If required, it can be brought into such a form by “undoing” the division of the first set of equations through $h^\alpha(\phi)$ and changing the signs in the last row.

This leads to the reduced system

$$\sum_{\alpha \in \mathcal{P}_p} \delta \tilde{c}^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^{N_p} \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \left(\mathbf{r}_{\tilde{\mu}^\alpha}^{(n)} + \delta \tilde{\mu} \right) h^\alpha(\phi) = \mathbf{r}_{\tilde{c}}$$

and thus, defining $\tilde{\zeta}^\alpha := \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \mathbf{r}_{\tilde{\mu}^\alpha}$, the $(K-1)$ -dimensional Schur complement system

$$\mathbf{S}_{\tilde{c}} \delta \tilde{\mu} = \mathbf{r}_{\tilde{c}}^{(n)} - \sum_{\alpha \in \mathcal{P}_p} \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \mathbf{r}_{\tilde{\mu}^\alpha}^{(n)} h^\alpha(\phi) = \mathbf{r}_{\tilde{c}}^{(n)} - \sum_{\alpha \in \mathcal{P}_p} (\tilde{\zeta}^\alpha)^{(n)} h^\alpha(\phi) \quad (7.36)$$

with the (reduced) Schur complement matrix

$$\mathbf{S}_{\tilde{c}} = \sum_{\alpha \in \mathcal{P}_p} \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} h^\alpha(\phi). \quad (7.37)$$

Once $\delta \tilde{\mu}$ (resp. $\tilde{\mu}^{(n+1)}$) is found, the new \tilde{c}^α -values are then easily recovered using Equation (7.35) in terms of the (independent) subsystems

$$(\tilde{c}^\alpha)^{(n+1)} = (\tilde{c}^\alpha)^{(n)} + \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \left(\mathbf{r}_{\tilde{\mu}^\alpha}^{(n)} + \delta \tilde{\mu} \right) = (\tilde{c}^\alpha)^{(n)} + \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \left(\tilde{\mu}^{(n+1)} - (\tilde{\mu}^\alpha)^{(n)} \right). \quad (7.38)$$

Remark 100. This procedure can in principle alternatively be interpreted as a predictor-corrector scheme. In fact, defining the auxiliary quantities

$$(\tilde{c}^\alpha)^{(n+\frac{1}{2})} := (\tilde{c}^\alpha)^{(n)} + \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \left(\tilde{\mu}^{(n)} - (\tilde{\mu}^\alpha)^{(n)} \right) \quad (7.39)$$

corresponding to the predicted values of the phase-specific concentration for the **current** value $\tilde{\mu}^{(n)}$ of $\tilde{\mu}$, the increment $\delta \tilde{\mu}$ and the final $(\tilde{c}^\alpha)^{(n+1)}$ satisfy

$$\mathbf{S}_{\tilde{c}} \delta \tilde{\mu} = \tilde{c} - \sum_{\alpha \in \mathcal{P}_p} (\tilde{c}^\alpha)^{(n+\frac{1}{2})} h^\alpha(\phi) = \mathbf{r}_{\tilde{c}}^{(n+\frac{1}{2})}$$

and

$$(\tilde{c}^\alpha)^{(n+1)} = (\tilde{c}^\alpha)^{(n+\frac{1}{2})} + \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \right)^{-1} \delta \tilde{\mu}.$$

It is easily verified that, after a division by $h^\alpha(\phi)$ for the phases with $h^\alpha(\phi) \neq 0$, this system is also the FONC of the minimization problem

$$\begin{cases} \text{minimize} & \frac{1}{2} \sum_{\alpha \in \mathcal{P}_p} \left((\tilde{c}^\alpha)^{(n+1)} - (\tilde{c}^\alpha)^{(n+\frac{1}{2})} \right) \cdot \frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2} \left((\tilde{c}^\alpha)^{(n)} \right) \cdot \left((\tilde{c}^\alpha)^{(n+1)} - (\tilde{c}^\alpha)^{(n+\frac{1}{2})} \right) h^\alpha(\phi) \\ \text{subject to} & \sum_{\alpha \in \mathcal{P}_p} (\tilde{c}^\alpha)^{(n+1)} h^\alpha(\phi) = \tilde{c} \end{cases}$$

and thus a correction step through a weighted (by the $\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{c}^\alpha)^2}$ -matrices evaluated for the old concentration values $(\tilde{c}^\alpha)^{(n)}$) projected operation onto the subspace of the phase-specific concentration values satisfying the sum-constraint $\sum_{\alpha \in \mathcal{P}_p} \tilde{c}^\alpha h^\alpha(\phi) = \tilde{c}$. \diamond

The block-factorization procedure using the Schur-complement $\mathbf{S}_{\tilde{\mathbf{c}}}$ instead of a direct solution of the full system at once has two advantages. On the one hand, it can take advantage of the block-diagonal structure. In particular, it is sufficient to work with the individual subblocks $\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\mathbf{c}}^\alpha)^2}$ alone when calculating the Schur complement $\mathbf{S}_{\tilde{\mathbf{c}}}$ and the modified residual in Equation (7.36) as well as for the $(\tilde{\mathbf{c}}^\alpha)^{(n+1)}$ in equation (7.38). On the other hand, while it is very expensive in terms of memory requirements to store all the phase-specific concentrations for the entire computational domain when a large number of phases are present, one can easily store the phase-independent multiplier $\tilde{\boldsymbol{\mu}}$. At any point at which it is necessary to use the \mathbf{c}^α , they can be recalculated from the **known** value $\tilde{\boldsymbol{\mu}}$, at least provided none of the other parameters (here ϕ and T) have changed in between. Starting e.g. from the initial guess $(\tilde{\mathbf{c}}^\alpha)^{(0)} = \mathbf{c}$, the corresponding Newton scheme reduces precisely to the same steps as the determination of the $(\mathbf{c}^\alpha)^{(n+\frac{1}{2})}$ in Equation (7.39) above, as one in fact has²⁷

$$(\tilde{\mathbf{c}}^\alpha)^{(n+1)} := (\tilde{\mathbf{c}}^\alpha)^{(n)} + \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\mathbf{c}}^\alpha)^2} \right)^{-1} \left(\tilde{\boldsymbol{\mu}} - (\tilde{\boldsymbol{\mu}}^\alpha)^{(n)} \right) = (\tilde{\mathbf{c}}^\alpha)^{(n)} + \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \left(\tilde{\boldsymbol{\mu}} - (\tilde{\boldsymbol{\mu}}^\alpha)^{(n)} \right). \quad (7.40)$$

One can therefore simply reuse the same functionality for both problems, that of determining the phase-specific concentrations and the reduced chemical potential from (ϕ, \mathbf{c}, T) alone or from (ϕ, \mathbf{c}, T) and $\tilde{\boldsymbol{\mu}}$, the latter one of course being somewhat cheaper as there is in particular no need of actually inverting the $\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\mathbf{c}}^\alpha)^2}$ -submatrices in order to determine $\mathbf{S}_{\tilde{\mathbf{c}}}$.

In matrix terms, the increments in the phase-specific concentration and $\tilde{\boldsymbol{\mu}}$ as a function of the residuals from Equation (7.34) can succinctly be summarized as²⁸ $\begin{pmatrix} (\delta \tilde{\mathbf{c}}^\alpha)_{\alpha \in \mathcal{P}_p} \\ \delta \boldsymbol{\mu} \end{pmatrix} = \mathbf{A} \begin{pmatrix} (\mathbf{r}_{\tilde{\boldsymbol{\mu}}^\alpha})_{\alpha \in \mathcal{P}_p} \\ \mathbf{r}_{\tilde{\mathbf{c}}} \end{pmatrix}$ with

$$\mathbf{A} = \begin{pmatrix} & & & & \begin{pmatrix} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^1}{\partial \tilde{\mathbf{c}}^1} \right)^{-1} \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \\ \left(\frac{\partial \tilde{\boldsymbol{\mu}}^2}{\partial \tilde{\mathbf{c}}^2} \right)^{-1} \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \\ \vdots \\ \left(\frac{\partial \tilde{\boldsymbol{\mu}}^{N_p}}{\partial \tilde{\mathbf{c}}^{N_p}} \right)^{-1} \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \end{pmatrix} \\ & \mathbf{A}_{11} & & & \\ \begin{pmatrix} -h^1 \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^1}{\partial \tilde{\mathbf{c}}^1} & -h^2 \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^2}{\partial \tilde{\mathbf{c}}^2} & \dots & -h^{N_p} \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^{N_p}}{\partial \tilde{\mathbf{c}}^{N_p}} \end{pmatrix} & & & \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \end{pmatrix}, \quad (7.41)$$

where the individual block-entries of the submatrix \mathbf{A}_{11} are given by

$$\mathbf{A}_{11}^{\alpha\beta} = \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \left(\delta^{\alpha\beta} \mathbf{I} - h^\beta(\phi) \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\beta}{\partial \tilde{\mathbf{c}}^\beta} \right)^{-1} \right). \quad (7.42)$$

While this representation can be convenient for theoretical purposes²⁹, its (1,1)-block - unlike the one corresponding to Equation (7.32) - is now a fully filled matrix due to the coupling by the

²⁷Where, as usual, it is preferable to solve these systems instead of multiplying the right-hand side by the inverse of $\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\mathbf{c}}^\alpha)^2} = \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}$.

²⁸This is simply a formula for the inverse of a matrix in terms of its Schur-complement in the case when the off-diagonal blocks are not transposes of each other, see e.g. [10] or [77].

²⁹Yet another representation of this inverse - which is closely linked with the predictor-corrector interpretation above - is based on the following block-factorization as a ‘‘perturbation’’ of the inverse of the (1,1)-block (see [77]),

$$\begin{pmatrix} \mathbf{A} & \mathbf{R} \\ \mathbf{L} & \mathbf{C} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \begin{pmatrix} -\mathbf{A}^{-1} \mathbf{R} \\ \mathbf{I} \end{pmatrix} \mathbf{S}^{-1} \begin{pmatrix} -\mathbf{L} \mathbf{A}^{-1} & \mathbf{I} \end{pmatrix},$$

which, in terms of the given quantities here, corresponds to

$$\begin{pmatrix} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \begin{pmatrix} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \\ \mathbf{I} \end{pmatrix} \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \begin{pmatrix} -h^\alpha \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} & \mathbf{I} \end{pmatrix}.$$

sum constraint. For practical purposes, it is therefore preferable to work with the “sequential” representation of this inverse consisting in first determining $\delta\tilde{\boldsymbol{\mu}}$ and then simply updating the $\tilde{\boldsymbol{c}}^\alpha$ through $\delta\tilde{\boldsymbol{c}}^\alpha = \left(\frac{\partial\tilde{\boldsymbol{\mu}}^\alpha}{\partial\tilde{\boldsymbol{c}}^\alpha}\right)^{-1} (\boldsymbol{r}_{\tilde{\boldsymbol{\mu}}^\alpha} + \delta\tilde{\boldsymbol{\mu}}^\alpha)$.

Remark 101. It should be noted that, due to the linearity in $\tilde{\boldsymbol{\mu}}$, disposing of a good “initial guess” $\tilde{\boldsymbol{\mu}}^{(0)}$ has no direct effect on the basic Newton scheme in the form of (7.32) as $\tilde{\boldsymbol{\mu}}^{(1)}$ is a function of the $(\tilde{\boldsymbol{c}}^\alpha)^{(0)}$ only. This is also easily seen based on equations (7.36) and (7.34), as the prefactor of $\tilde{\boldsymbol{\mu}}^{(n)}$ in the update formula (7.36) is (due to the independence of $\tilde{\boldsymbol{\mu}}$ on α) precisely $-\boldsymbol{S}_{\tilde{\boldsymbol{c}}}$, i.e. one can equivalently replace Equation (7.36) by the alternative update rule

$$\boldsymbol{S}_{\tilde{\boldsymbol{c}}}(\tilde{\boldsymbol{\mu}}^{(n)} + \delta\tilde{\boldsymbol{\mu}}) = \boldsymbol{S}_{\tilde{\boldsymbol{c}}}\tilde{\boldsymbol{\mu}}^{(n+1)} = (\tilde{\boldsymbol{c}} - \sum_{\alpha \in \mathcal{P}_p} \tilde{\boldsymbol{c}}^\alpha h^\alpha(\boldsymbol{\phi})) + \sum_{\alpha \in \mathcal{P}_p} \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\boldsymbol{c}}^\alpha)^2} \right)^{-1} (\tilde{\boldsymbol{\mu}}^\alpha)^{(n)} h^\alpha(\boldsymbol{\phi}).$$

With $\tilde{\boldsymbol{\mu}}^{(n+1)}$, the $(\delta\tilde{\boldsymbol{c}}^\alpha)_{\alpha \in \mathcal{P}_p}$ can then be recovered using the left expression in Equation (7.35). \diamond

Remark 102. A natural initial guess is to simply choose the initial phase-specific concentrations to be the same as the total concentration³⁰, i.e. setting $(\tilde{\boldsymbol{c}}^\alpha)^{(0)} = \tilde{\boldsymbol{c}}$, $\alpha = 1, \dots, \mathcal{P}_p$. Due to the averaging property $\sum_{\alpha \in \mathcal{P}_p} \tilde{\boldsymbol{c}}^\alpha h^\alpha(\boldsymbol{\phi}) = (\sum_{\alpha \in \mathcal{P}_p} h^\alpha(\boldsymbol{\phi}))\tilde{\boldsymbol{c}} = \tilde{\boldsymbol{c}}$, this in particular ensures that the sum-constraint on the concentrations is automatically satisfied for the initial guess, i.e. $\boldsymbol{r}_{\tilde{\boldsymbol{c}}}^{(0)} = \mathbf{0}$.

Due to the linearity of the second equation in either Equation (7.33) or (7.32) w.r.t. the $\delta\boldsymbol{c}^\alpha$ and a simple recursion, this property will in fact be inherited by all subsequent iterates, as, per construction, $\sum_{\alpha \in \mathcal{P}_p} \delta\tilde{\boldsymbol{c}}^\alpha h^\alpha(\boldsymbol{\phi}) = \tilde{\boldsymbol{c}} - \sum_{\alpha \in \mathcal{P}_p} (\tilde{\boldsymbol{c}}^\alpha)^{(n)} h^\alpha(\boldsymbol{\phi})$ and the weighted average of the \boldsymbol{c}^α does therefore not change between the n -th and $(n+1)$ -th step provided it was already consistent³¹. This unfortunately does not allow for any major simplifications of the algorithm above, as, despite the ability of dropping the right-most row in Equation (7.41), the (larger) $(1, 1)$ -block given by \boldsymbol{A}_{11} in Equation (7.42) still contains all inverse matrices appearing in (7.41). \diamond

Remark 103. Complementing Remark 101, if one believes an old value of $\tilde{\boldsymbol{\mu}}$ to be a good approximation of the new one - as is e.g. the case in many time-stepping schemes, in particular if the time-step is small - one may nevertheless use this value to generate a well-educated initial guess for the $\tilde{\boldsymbol{c}}^\alpha$. To do so, it suffices to (approximately) determining the \boldsymbol{c}^α such that $\frac{\partial \tilde{f}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha} = \tilde{\boldsymbol{\mu}}$ which, as $\tilde{\boldsymbol{\mu}}$ is taken as fixed, only involves the small block-diagonal submatrices $\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha}$ whose factorizations need to be determined anyway.

In contrast to Remark 102, this does not ensure the fulfillment of the averaging on the $(\tilde{\boldsymbol{c}}^\alpha)^{(0)}$ to $\tilde{\boldsymbol{c}}$, but, if solved exactly³², would ensure the equality of the $\tilde{\boldsymbol{\mu}}^\alpha$ with $\tilde{\boldsymbol{\mu}}$ and thus the vanishing of $(\boldsymbol{r}_{\tilde{\boldsymbol{\mu}}^\alpha}^{(0)})_{\alpha \in \mathcal{P}_p}$.

While the contribution of the first (larger) part of \boldsymbol{A} to the increments does in this case vanish in the first step - i.e. it seems as if one could make full use of the sparsity of the $(1, 1)$ -submatrix $\left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\boldsymbol{c}}^\beta}\right)_{\alpha, \beta \in \mathcal{P}_p}$ in the first equation of the system (7.32) instead of the fully filled \boldsymbol{A}_{11} -matrix in Equation (7.42) - this, similar to Remark 102, typically would generally not lead to any substantial decrease in the computational effort in the first iteration step as compared to solving the full system in Equation (7.32). In fact, on the one hand, the appearance of $\boldsymbol{S}_{\tilde{\boldsymbol{c}}}$ still

³⁰Which, as the conserved quantity, is a priori the natural one to be stored. In relation with the discussion in Section 7.1.4, an alternative consists in - as proposed e.g. in [19] - instead storing $\tilde{\boldsymbol{\mu}}$ only and using the (discretized) evolution equation for $\tilde{\boldsymbol{\mu}}$ from Equation (7.24) for updating $\tilde{\boldsymbol{\mu}}$ instead of the concentration itself. While this approach is in certain cases quite efficient and has been used quite successfully also for very large simulation setups (see for example the works [36], [70], [39] and the references therein), it is also somewhat restrictive as will be discussed below.

³¹The same will actually happen for any initial guess after the first iteration.

³²This is certainly not to be recommended in general! Even if the individual sub-iterations in this first step would be slightly less inexpensive (but see Remark 103 below) than an iteration on the full system, this subiteration on a **subsystem** of the actual one to be solved would be at the cost of the quadratic convergence of the Newton-scheme on the **actual** system to be solved once the $\tilde{\boldsymbol{c}}^\alpha$ and $\tilde{\boldsymbol{\mu}}$ are sufficiently close to the true ones.

requires an factorization and inversion of the submatrices $\left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha}\right)_{\alpha \in \mathcal{P}_p}$ for all phases in \mathcal{P}_p . On the other hand, the natural algorithmic approach for the solution of the full system consists in **first** determining $\delta \tilde{\boldsymbol{\mu}}$ in terms of the auxiliary vectors $\tilde{\boldsymbol{\zeta}}^\alpha$ as in Equation (7.36) and then using Equation (7.38) for determining the $\delta \tilde{\boldsymbol{c}}^\alpha$. This requires precisely the same steps as the first iteration step assuming, as above, that the $\boldsymbol{r}_{\tilde{\boldsymbol{\mu}}^\alpha}$, $\alpha \in \mathcal{P}_p$ vanish, except for the solution of Equation (7.38) based on the already factorized or inverted submatrices $\left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha}\right)_{\alpha \in \mathcal{P}_p}$. \diamond

Remark 104. In relation with the approaches in [19] and [56], it should be kept in mind that there is a seemingly very convenient alternative solution procedure, which consists in directly choosing the (common) multiplier $\tilde{\boldsymbol{\mu}}$ as the primary unknown and expressing the remaining unknowns $\tilde{\boldsymbol{c}}^\alpha$ as direct functions of $\tilde{\boldsymbol{\mu}}$ and T .

This reduction in the number of unknowns is clearly a quite appealing feature, since it replaces the system (7.10) with one consisting of the sum-constraint as a function of $\tilde{\boldsymbol{\mu}}$ alone,

$$\sum_{\alpha=1}^N \tilde{\boldsymbol{c}}^\alpha(\tilde{\boldsymbol{\mu}}, T) h^\alpha(\boldsymbol{\phi}) \stackrel{!}{=} \tilde{\boldsymbol{c}}. \quad (7.43)$$

Furthermore, taking again a Newton scheme as the obvious choice of solving Equation (7.43), one directly obtains the update rule

$$\left(\sum_{\alpha=1}^N \frac{\partial \tilde{\boldsymbol{c}}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}(\tilde{\boldsymbol{\mu}}^{(n)}, T) h^\alpha(\boldsymbol{\phi})\right) \delta \tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{c}} - \sum_{\alpha=1}^N \tilde{\boldsymbol{c}}^\alpha(\tilde{\boldsymbol{\mu}}^{(n)}, T) h^\alpha(\boldsymbol{\phi}). \quad (7.44)$$

This can in fact be a very efficient approach provided one disposes of a direct formulation of the $\tilde{\boldsymbol{c}}^\alpha$ as functions of $\tilde{\boldsymbol{\mu}}$ and T since it only requires the solution of a single and usually fairly small system for the $K - 1$ unknowns in $\delta \tilde{\boldsymbol{\mu}}$.

It is much less practical though if this is not the case, i.e. if one is not able to derive an explicit expression for $\tilde{\boldsymbol{c}}^\alpha(\tilde{\boldsymbol{\mu}}, T)$. Even though this does not preclude the use of such a scheme since one can always solve the equation

$$\frac{\partial \tilde{f}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha}(\tilde{\boldsymbol{c}}^\alpha, T) = \tilde{\boldsymbol{\mu}} \quad (7.45)$$

numerically³³ and then recover the required expressions for $\frac{\partial \tilde{\boldsymbol{c}}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}$ by differentiating Equation (7.45) for $\boldsymbol{c}^\alpha = \boldsymbol{c}^\alpha(\tilde{\boldsymbol{\mu}}, T)$ leading to

$$\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{\boldsymbol{c}}^\alpha)^2}(\tilde{\boldsymbol{c}}^\alpha, T) \frac{\partial \boldsymbol{c}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}(\tilde{\boldsymbol{\mu}}, T) = \boldsymbol{I} \quad \text{resp.} \quad \frac{\partial \boldsymbol{c}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}(\tilde{\boldsymbol{\mu}}, T) = \left(\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{\boldsymbol{c}}^\alpha)^2}(\tilde{\boldsymbol{c}}^\alpha(\tilde{\boldsymbol{\mu}}, T), T)\right)^{-1}, \quad (7.46)$$

this reintroduces the necessity of the inversion of these phase-specific matrices and leads back to the formula for the reduced Schur-complement matrix $\boldsymbol{S}_{\tilde{\boldsymbol{c}}}$ in Equation (7.37). Furthermore, the numerical solution of Equation (7.45) is again most naturally done using a Newton-scheme and thus through the update rule $(\tilde{\boldsymbol{c}}^\alpha)^{(n+1)} = (\tilde{\boldsymbol{c}}^\alpha)^{(n)} + \delta \tilde{\boldsymbol{c}}^\alpha$ with $\delta \tilde{\boldsymbol{c}}^\alpha$ determined by

$$\frac{\partial^2 \tilde{f}^\alpha}{\partial (\tilde{\boldsymbol{c}}^\alpha)^2}((\tilde{\boldsymbol{c}}^\alpha)^{(n)}, T) \delta \tilde{\boldsymbol{c}}^\alpha = \tilde{\boldsymbol{\mu}} - \frac{\partial \tilde{f}^\alpha}{\partial \tilde{\boldsymbol{c}}^\alpha}((\tilde{\boldsymbol{c}}^\alpha)^{(n)}, T), \quad (7.47)$$

reintroducing the \boldsymbol{c}^α and $\delta \boldsymbol{c}^\alpha$ as auxiliary unknowns. In addition, Equation (7.45) should a priori be solved exactly in order for the expression for $\frac{\partial \tilde{\boldsymbol{c}}^\alpha}{\partial \tilde{\boldsymbol{\mu}}}$ in Equation (7.46) to be correct, and will therefore in general require several Newton steps.

³³This is also indicated in [56].

The additional computational cost of this can of course be alleviated by deliberately allowing for inaccuracies in the solution of this system by e.g. replacing the exact solution of the equations (7.45) with a single Newton-step as in Equation (7.47), but one has to exercise some care when doing so. While simply using the relation (7.46) is in principle still possible, this will seriously impede the convergence of the update rule (7.44). The problem with this approach is that it neglects an effect distinguishing Equation (7.47) from the derivation of Equation (7.46). Namely, whereas Equation (7.46) is based on predicting the change of \tilde{c}^α to changes in $\tilde{\mu}$ $\frac{\partial^2 f^\alpha}{\partial(\tilde{c}^\alpha)^2} \delta\tilde{c}^\alpha = \delta\tilde{\mu}$, what actually happens due the update (7.47) is that $\delta\tilde{c}^\alpha$ changes according to

$$\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{c}^\alpha)^2} \left((\tilde{c}^\alpha)^{(n)}, T \right) \delta\tilde{c}^\alpha = \delta\tilde{\mu} + \tilde{\mu} - \frac{\partial \tilde{f}^\alpha}{\partial \tilde{c}^\alpha} \left((\tilde{c}^\alpha)^{(n)}, T \right),$$

and therefore is not only affected by $\delta\tilde{\mu}$ but also by the residual $r_{\tilde{\mu}^\alpha}$ for Equation (7.45) at the last step. Including this effect for restoring the quadratic convergence rate then simply leads back to Equation (7.36) and thus is ultimately a somewhat cumbersome and more error-prone way of obtaining the the previous linearized system (7.32) obtained through a “mechanical” differentiation based on the direct use of the more natural (for the \tilde{f}^α given as functions of \tilde{c}^α and T) variables. \diamond

Remark 105. By the previous remark, the computational cost of calculating an increment in $\tilde{\mu}$ due to a change in the concentrations using implicitly defined functions \tilde{c}^α in terms of $\tilde{\mu}$ and T is, at least when using the inaccurate version, roughly similar to the one of obtaining $\tilde{\mu}$ based on a given concentration as in Section 7.1.5 based on Equation (7.10).

A very important exception to this occurs in the bulk-regions, where, given $\tilde{f}^\alpha(\tilde{c}^\alpha, T)$ and the average concentration \tilde{c} and the temperature T , one can obtain $\tilde{\mu}$ directly from $\tilde{\mu} = \frac{\partial \tilde{f}^\alpha}{\partial \tilde{c}^\alpha}(\tilde{c}, T)$ and thus is an operation which is expected to be relatively cheap. This is a major advantage over a scheme (such as proposed in [19]) based on the storage of $\tilde{\mu}$ only, unless of course one has an explicit formulation of the \tilde{c}^α as functions of $\tilde{\mu}$ allowing to make use of the same simplification through a then also relatively cheap “conversion” of $\tilde{\mu}$ to \tilde{c}^α for the bulk-phase. If such is not the case, the storage of $\tilde{\mu}$ instead of \tilde{c} is unlikely to be an efficient approach as this enforces an iterative update scheme in some form for updating $\tilde{\mu}$ if there are any changes in the concentration. In particular, there is no really valid way of justifying the use of a “blind” single-step scheme as for the time-discretization of Equation (7.24) in [19]. Even if the changes in $\tilde{\mu}$ are very small, due to e.g. a small time-step, and a single Newton-step for $\delta\tilde{\mu}$ is therefore expected to be very accurate, this requires at least an accurate estimate of $\frac{\partial \tilde{c}^\alpha}{\partial \tilde{\mu}}$, and therefore a sufficiently accurate estimate of \tilde{c}^α in the bulk-phase for which there is none if the concentration is not stored³⁴.

Except for sufficiently simple free energies where one can derive the appropriate relations analytically, a formulation directly in terms of \tilde{c} or, if using a primarily $\tilde{\mu}$ -based formulation, at least the additional storage of the concentration field is therefore to be recommended. \diamond

³⁴One of course has to be careful about the meaning of sufficiently accurate and to take the changes of other parameters into account in order to maintain a high accuracy with such a scheme.

The Non-Reduced Formulation

Solving for the unknowns in the non-reduced formulation from Equation (7.9) is a priori a little more difficult. In this case, the linearized system is given by

$$\frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}((\mathbf{c}^\alpha)^{(n)}, T) + \frac{\partial^2 f^\alpha}{\partial (\mathbf{c}^\alpha)^2}((\mathbf{c}^\alpha)^{(n)}, T) \delta \mathbf{c}^\alpha - (\boldsymbol{\mu}^{(n)} + \delta \boldsymbol{\mu} + (\lambda^\alpha)^{(n+1)} \mathbf{e}) = \mathbf{0} \quad (7.48)$$

together with the sum-constraint $\sum_{\alpha \in \mathcal{P}_p} \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c}$ and the sum-constraints $(\mathbf{e} \cdot \mathbf{c}^\alpha)^{(n+1)} = 1$ on the phase-specific concentrations. On the one hand, this system a priori has the disadvantage of being somewhat larger due to the presence of the additional phase-specific multipliers $\boldsymbol{\lambda}$ and the additional sum-constraints on the \mathbf{c}^α . On the other hand, as pointed out above, due to the partial redundancy in the sum-constraint for \mathbf{c} , one needs to impose an additional condition on $\boldsymbol{\mu}$ in order to obtain unicity of the multipliers.

Nevertheless, this formulation is also somewhat more natural in the sense that the phase-specific free energies are often more easily described using the full concentration vectors \mathbf{c}^α instead of an ‘‘artificial’’ reduction to $K - 1$ components only.

Fortunately, this is not a very serious concern as one already disposes of a simple solution procedure for the reduced formulation and the system in Equation (7.48) together with the constraints provides all the differential information on the system required to recover the remaining quantities if desired. Both the reduced and the non-reduced case can therefore be handled in a relatively transparent fashion using the same central algorithm.

More precisely, subtracting the K 'th row in the System (7.48) from the first $K - 1$ ones - corresponding to multiplying the system from the left by the matrix $\begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix}$ - makes it possible to eliminate the common factor λ^α and reduces the remaining system to

$$\begin{aligned} & \left(\frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_i^\alpha} \right)_{1 \leq i \leq K-1} - \frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_K^\alpha} \tilde{\mathbf{e}} + \begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix} \frac{\partial^2 f^\alpha}{\partial (\mathbf{c}^\alpha)^2}((\mathbf{c}^\alpha)^{(n)}, T) \delta \mathbf{c}^\alpha \\ & - \begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix} (\boldsymbol{\mu}^{(n)} + \delta \boldsymbol{\mu}) = \mathbf{0}. \end{aligned}$$

Firstly, $\left(\frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_i^\alpha} \right)_{1 \leq i \leq K-1} - \frac{\partial f^\alpha(\mathbf{c}^\alpha, T)}{\partial c_K^\alpha} \tilde{\mathbf{e}} = \frac{\partial f^\alpha((\mathbf{c}^\alpha)^{(n)}, 1 - \sum_{i=1}^{K-1} (c_i^\alpha)^{(n)}, T)}{\partial \tilde{\mathbf{c}}^\alpha} = \tilde{\boldsymbol{\mu}}^\alpha((\tilde{\mathbf{c}}^\alpha)^{(n)}, T)$ by Equation (7.11) and similarly $\mu_i - \mu_K = \tilde{\mu}_i$ and $\delta \mu_i - \delta \mu_K = \delta \tilde{\mu}_i$. Secondly, making use of the constraint to express δc_K^α as $-\sum_{i=1}^{K-1} \delta c_i^\alpha$ - corresponding to $\delta \mathbf{c}^\alpha = \begin{pmatrix} \mathbf{I}_{K-1} \\ -\tilde{\mathbf{e}} \end{pmatrix} \delta \tilde{\mathbf{c}}^\alpha$ - the system is reduced to the $(K - 1) \times (K - 1)$ one

$$\tilde{\boldsymbol{\mu}}^\alpha((\tilde{\mathbf{c}}^\alpha)^{(n)}, T) + \begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix} \frac{\partial^2 f^\alpha}{\partial (\mathbf{c}^\alpha)^2}((\mathbf{c}^\alpha)^{(n)}, T) \begin{pmatrix} \mathbf{I}_{K-1} \\ -\tilde{\mathbf{e}} \end{pmatrix} \delta \tilde{\mathbf{c}}^\alpha - (\tilde{\boldsymbol{\mu}}^{(n)} + \delta \tilde{\boldsymbol{\mu}}) = \mathbf{0}. \quad (7.49)$$

This, as is to be expected, is precisely the one in Equation (7.32) that would be obtained based on directly using a reduced formulation, except that $\tilde{\boldsymbol{\mu}}^\alpha((\tilde{\mathbf{c}}^\alpha)^{(n)}, T)$ and the matrix $\frac{\partial \tilde{\boldsymbol{\mu}}}{\partial \tilde{\mathbf{c}}^\alpha}$ is recovered a posteriori as

$$\tilde{\boldsymbol{\mu}}^\alpha((\tilde{\mathbf{c}}^\alpha)^{(n)}, T) = \begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix} \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}((\mathbf{c}^\alpha)^{(n)}, T) \quad (7.50)$$

and

$$\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}((\tilde{\mathbf{c}}^\alpha)^{(n)}, T) = \begin{pmatrix} \mathbf{I}_{K-1} & -\tilde{\mathbf{e}} \end{pmatrix} \frac{\partial^2 f^\alpha}{\partial (\mathbf{c}^\alpha)^2}((\mathbf{c}^\alpha)^{(n)}, T) \begin{pmatrix} \mathbf{I}_{K-1} \\ -\tilde{\mathbf{e}} \end{pmatrix} \quad (7.51)$$

through two simple algebraic operations instead of a priori being included in the formulation of the f^α .

Remark 106. Even though it may seem relatively cumbersome to algebraically reproduce a calculation which could from the outset be done by hand when constructing the free energy function, this is also closely related to a matter of “convention”. In particular, it is clear that starting from a reduced formulation is in many ways the simplest approach, in particular since one typically does not face any of the issues associated with such a formulation in the obstacle case for the phasefield (i.e. where phases may not be allowed to move according to the changes of the other ones). Nevertheless, this also introduces a certain asymmetry into the formulation by having to choose which component is the one being eliminated (the K -th one clearly being a good candidate though simply in terms of memory management and simplicity of any loops). This may on the one hand be perceived as unpleasant on aesthetic grounds and on the other hand entails a strict consistency requirement.

One advantage of the purely algebraic procedure above is its high robustness with respect to these issues, since, even though it entails an in principle completely unnecessary calculation if used with a reduced formulation (the operations in Equation (7.50) and (7.51) are then easily seen to have no effect), it allows combining various preexisting model components even if based on different conventions³⁵ \diamond

In relation with the non-reduced formulation and the reduced nature of the Newton-steps in Equation (7.49), there are now essentially two manners to proceed. If one is solely interested in an evaluation of the final concentrations, either for their own values or for evaluating the f -function itself, it is obviously sufficient to simply recover the value of the last component c_K^α in the solution of the reduced system through the sum-constraint as $1 - \sum_{i=1}^{K-1} c_i^\alpha$.

If one is instead also interested in the values of $\boldsymbol{\mu}$ and the λ^α , one can also update these after the last step of the solution procedure³⁶. More precisely, as the solution satisfies (up to the chosen tolerance for the residual of course) $\boldsymbol{\mu}^\alpha = \boldsymbol{\mu} + \lambda^\alpha \mathbf{e}$, one can easily recover the required information from the final values of $\boldsymbol{\mu}^\alpha$ at which the solution was considered sufficiently accurate for the three choices of restriction on $\boldsymbol{\mu}$ outlined in Subsection 7.1.2 since this relation implies $\mathbf{e} \cdot \boldsymbol{\mu}^\alpha = \mathbf{e} \cdot \boldsymbol{\mu} + K\lambda^\alpha$ with a known left-hand side. If $\boldsymbol{\mu}$ is fixed by enforcing its average to equal 0, the second term drops out and one has $\lambda^\alpha = \frac{1}{K} \mathbf{e} \cdot \boldsymbol{\mu}^\alpha$, i.e. λ^α is simply equal to the average of the $\boldsymbol{\mu}^\alpha$. In contrast, if $\boldsymbol{\mu}$ is fixed by the condition that $\mu_K = 0$, one obviously has $\tilde{\boldsymbol{\mu}} = (\mu_i)_{1 \leq i \leq K-1} - \mu_K \tilde{\mathbf{e}} = (\mu_i)_{1 \leq i \leq K-1}$, and thus $\tilde{\mathbf{e}} \cdot \tilde{\boldsymbol{\mu}} = \mathbf{e} \cdot \boldsymbol{\mu}$ as the last entry vanishes, therefore leading to $\lambda^\alpha = \frac{1}{K} (\mathbf{e} \cdot \boldsymbol{\mu}^\alpha - \tilde{\mathbf{e}} \cdot \tilde{\boldsymbol{\mu}})$. Finally, if $\boldsymbol{\mu}$ is chosen indirectly as the h^α -weighted average of the $\boldsymbol{\mu}^\alpha = \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}$, there is no need to first determine the λ^α as $\boldsymbol{\mu}$ can be obtained directly from the $\boldsymbol{\mu}^\alpha$. The λ^α then immediately follow from the equality $\lambda^\alpha = \frac{1}{K} \mathbf{e} \cdot (\boldsymbol{\mu}^\alpha - \boldsymbol{\mu})$.

Remark 107. Note that despite these different choices, this has no effect on the final concentration values, and that it does also not affect the evolution of the concentration by the construction of the L_{ij} -matrix in [52]. \diamond

³⁵For example, it is possible to combine a chemical contribution formulated in a reduced fashion through the first $K - 1$ components with other contributions where the only direct contribution arises through the K -th component. This of course has no theoretical relevance, but avoids some very tedious issues at a practical level.

³⁶As the basic Newton-scheme is shifted to the reduced formulation, this is typically also where the residual would typically be controlled in its reduced form, i.e. one never directly uses $\boldsymbol{\mu}$ and the λ^α in the algorithm itself.

7.1.6 An Extension to the Non-Isothermal Case

The “grand-chemical” model considered in the previous section seems to have previously only been used in either an isothermal setting, or, under a relatively strong simplifying assumption, in combination with a prescribed non-uniform temperature field. In the form above, it is relatively clear though how the more quantitative approach can be formulated in a non-isothermal setting. With the maximization of the entropy functional from Equation (3.3) replacing the minimization of the free energy functional in Equation (6.10) as the appropriate variational principle, it is natural to define the entropy density within an interface region by

$$s(\phi, \mathbf{c}, e) := \max_{(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}, e)} \left\{ \sum_{\alpha=1}^N s^\alpha(\mathbf{c}^\alpha, e^\alpha) h^\alpha(\phi) \right\}, \quad (7.52)$$

using the phase-specific quantities $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ instead of the average quantities (\mathbf{c}, e) and with the admissible set $\mathcal{A}(\phi, \mathbf{c}, e)$ consisting of all $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ such that³⁷

$$\begin{cases} \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c}, \\ \sum_{i=1}^K c_i^\alpha = 1 \quad \forall \alpha, \\ \sum_{\alpha=1}^N e^\alpha h^\alpha(\phi) = e \end{cases} \quad (7.53)$$

replacing the one in Equation (7.4).

The corresponding optimality conditions can again be derived based on the modified Lagrange-function

$$\begin{aligned} L(\phi, \mathbf{c}, e, (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, (e^\alpha)_{1 \leq \alpha \leq N}, \boldsymbol{\eta}, \hat{\boldsymbol{\lambda}}, \beta) &= \sum_{\alpha=1}^N s^\alpha(\mathbf{c}^\alpha, e^\alpha) h^\alpha(\phi) - \beta \cdot \left(\sum_{\alpha=1}^N h^\alpha \mathbf{c}^\alpha - \mathbf{c} \right) \\ &\quad - \sum_{\alpha=1}^N \hat{\lambda}^\alpha \left(\sum_{i=1}^K c_i^\alpha - 1 \right) - \beta \left(\sum_{\alpha=1}^N h^\alpha e^\alpha - e \right) \end{aligned}$$

with $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^N$ and (a priori) $\beta \in \mathbf{R}$. The first-order necessary conditions can then be obtained from $\frac{\partial L}{\partial \mathbf{c}^\alpha} = \mathbf{0}$ and $\frac{\partial L}{\partial e^\alpha} = 0$, leading, after by a division through h^α for all $\alpha \in \mathcal{P}_p$, to the local quasi-equilibrium conditions

$$\begin{cases} \frac{\partial s^\alpha}{\partial \mathbf{c}^\alpha} = \boldsymbol{\eta} + \hat{\lambda}^\alpha \mathbf{e}, \\ \frac{\partial s^\alpha}{\partial e^\alpha} = \beta. \end{cases} \quad (7.54)$$

As $\frac{\partial s^\alpha}{\partial e^\alpha} = \frac{1}{T^\alpha}$ and in analogy to the equilibrium condition $T = \text{const}$ for the temperature field, it is of course not surprising that the second condition implies the equality of the phase-specific temperatures $T^\alpha = \frac{1}{\beta} =: T$ for all phases (or, more precisely, for those with $h^\alpha(\phi) \neq 0$). Based upon the monotonicity condition $\frac{\partial s^\alpha}{\partial e^\alpha} > 0$ for the bulk entropy densities, this common value in addition is necessarily positive, and one can substitute $\boldsymbol{\eta}$ and $\hat{\lambda}^\alpha$ by $-\frac{\boldsymbol{\mu}}{T^\alpha} = -\frac{\boldsymbol{\mu}}{T}$ and $-\frac{\lambda^\alpha}{T^\alpha} = -\frac{\lambda^\alpha}{T}$ and cancel the T in the first equation above to obtain the local quasi-equilibrium conditions in the more “pleasant” form

$$T^\alpha = T \quad \text{and} \quad \boldsymbol{\mu}^\alpha = \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, T) = \boldsymbol{\mu} + \lambda^\alpha \mathbf{e} \quad \forall \alpha \in \mathcal{P}_p \quad (7.55)$$

and thus in particular including the previous condition on the chemical potentials as in the isothermal case.

Remark 108. The equilibrium conditions above are clearly what one would expect in relation with the previous discussion and based on the close analogy between the optimization problems

³⁷Note that it is again assumed here that the positivity constraints on the \mathbf{c}^α can be safely neglected.

in Equations (7.52) and (7.3). On the one hand, as discussed in e.g. [56], it is the compatibility of the (local) assumption $T^\alpha = T$ with the (global) equilibrium condition $T = \text{const}$ which can be used to explain the much larger success of the early phasefield models for purely temperature-driven solidification problems as compared to concentration-driven ones. On the other hand, the equilibration of the chemical potentials is well-known to be the relevant equilibrium condition for the concentration field, whether or not the process under consideration is from the outset assumed to be isothermal.

The point of the (re-)derivation of these conditions based upon the definition (7.52) is therefore not to show that these conditions are reasonable ones to impose within the interface region, but rather to verify that they arise naturally through a definition of the relevant local density which is compatible with the underlying global variational principle. \diamond

Remark 109. It should be noted that T^α and μ^α in Equation (7.55) are a priori - based upon their derivation from (7.54) - to be considered as functions of \mathbf{c}^α and e^α (and in particular μ^α as $\mu^\alpha(\mathbf{c}^\alpha, e^\alpha) = -T^\alpha(\mathbf{c}^\alpha, e^\alpha) \frac{\partial s^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha, e^\alpha)$). As it is often preferred in practice to work with more concrete temperature values T^α instead of the e^α as the “independent” variable, this is not necessarily the most convenient form and one may choose to work, using the ubiquitous changes of variables in the thermodynamic setting, with a different set of “primary” unknowns. This does of course not change the content of the equilibrium conditions, but, as it changes their representation, can affect the solution process. Some advantages and disadvantages of such changes of variables will be discussed after taking a closer look at the link of the Definition (7.52) of the entropy $s(\phi, \mathbf{c}, e)$ with the Definition (7.3) from the previous section. \diamond

Some Basic Thermodynamic Relations

Before briefly discussing the solution procedure for the local system (7.55) in terms of the (assumed to be given) phasefield ϕ and the conserved phase-averaged quantities \mathbf{c} and e , it is helpful to recover a number of basic thermodynamic properties resembling those from the more standard case which are in fact implied by the Definition (7.52).

The primary difficulty here, as in the isothermal case, lies in the fact that the phase-inherent quantities $(\mathbf{c}^\alpha, e^\alpha)$ upon which $s(\phi, \mathbf{c}, e)$ is based are now only defined implicitly as the ones fulfilling the local quasi-equilibrium conditions (7.55) above, i.e. $(\mathbf{c}^\alpha, e^\alpha) = (\mathbf{c}^\alpha(\phi, \mathbf{c}, e), e^\alpha(\phi, \mathbf{c}, e))$. In particular due to the equilibrium condition on the T^α , a first useful step is to obtain the corresponding free energy density $f(\phi, \mathbf{c}, T)$, defined by

$$f(\phi, \mathbf{c}, T) = \inf_e \{e - Ts(\phi, \mathbf{c}, e)\}. \quad (7.56)$$

It turns out that the following intuitively pleasing but not entirely obvious characterization holds for $f(\phi, \mathbf{c}, T)$:

Lemma 7. *Let $s(\phi, \mathbf{c}, e)$ be defined by Equation (7.52) and assume that the local quasi-equilibrium conditions (7.55) allow for a unique (up to the addition of an arbitrary constant vector $\xi \mathbf{e}$ to $\boldsymbol{\mu}$) solution which depends smoothly on \mathbf{c} and e and that the bulk entropy densities depend smoothly on \mathbf{c}^α and e^α . Then:*

- For any minimizer e in Equation (7.56), the common value of the phase-specific temperatures T^α for the solution of the local quasi-equilibrium conditions (7.55) coincide with the “external” parameter T .
- The free energy density defined in Equation (7.56) coincides with the one defined in Equation (7.3)

$$f(\phi, \mathbf{c}, T) = \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})} \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi). \quad (7.57)$$

Proof. • Firstly, under the smoothness assumption of the lemma,

$$s(\phi, \mathbf{c}, \mathbf{e}) = \max_{(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}, \mathbf{e})} \left\{ \sum_{\alpha} s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) h^{\alpha}(\phi) \right\} = \sum_{\alpha=1}^N s^{\alpha}(\mathbf{c}^{\alpha}(\phi, \mathbf{c}, \mathbf{e}), T^{\alpha}(\phi, \mathbf{c}, \mathbf{e})) h^{\alpha}(\phi)$$

depends smoothly on \mathbf{e} and thus any minimizer \mathbf{e} in definition (7.56) has to satisfy

$$\frac{1}{T} = \frac{\partial s}{\partial \mathbf{e}} = \sum_{\alpha} \left(\frac{\partial s^{\alpha}}{\partial \mathbf{c}^{\alpha}} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{e}} + \frac{\partial s^{\alpha}}{\partial e^{\alpha}} \frac{\partial e^{\alpha}}{\partial \mathbf{e}} \right) h^{\alpha}(\phi) = \sum_{\alpha} \left(-\frac{\boldsymbol{\mu}^{\alpha}}{T^{\alpha}} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{e}} + \frac{1}{T^{\alpha}} \frac{\partial e^{\alpha}}{\partial \mathbf{e}} \right) h^{\alpha}(\phi).$$

By the equilibrium conditions, T^{α} and $\frac{\boldsymbol{\mu}^{\alpha}}{T^{\alpha}}$ can be replaced by a common value θ (which at this point is not yet known to be the same as T) and $\frac{\boldsymbol{\mu} + \lambda^{\alpha} \mathbf{e}}{\theta}$, leading to

$$\frac{1}{T} = \sum_{\alpha} \left(-\frac{\boldsymbol{\mu} + \lambda^{\alpha} \mathbf{e}}{\theta} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{e}} + \frac{1}{\theta} \frac{\partial e^{\alpha}}{\partial \mathbf{e}} \right) h^{\alpha}(\phi).$$

As θ does not depend on α and ϕ does not depend on \mathbf{e} , the summation over the last term can be simplified to

$$\frac{1}{\theta} \frac{\partial \left(\sum_{\alpha=1}^N e^{\alpha} h^{\alpha}(\phi) \right)}{\partial \mathbf{e}} = \frac{1}{\theta} \frac{\partial e}{\partial \mathbf{e}} = \frac{1}{\theta}$$

and one in particular does not require to know $\frac{\partial e^{\alpha}}{\partial \mathbf{e}}$ anymore. In a similar manner, the summation over the terms $\frac{\boldsymbol{\mu}}{\theta} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{e}}$ leads to $\frac{\boldsymbol{\mu}}{\theta} \cdot \frac{\partial \left(\sum_{\alpha=1}^N \mathbf{c}^{\alpha} h^{\alpha}(\phi) \right)}{\partial \mathbf{e}} = \frac{\boldsymbol{\mu}}{\theta} \cdot \frac{\partial \mathbf{c}}{\partial \mathbf{e}} = 0$ and therefore simply drops out. That the same actually also happens for the (phase-dependent) term in λ^{α} is directly based upon the sum-constraint $\sum_{i=1}^K \mathbf{c}^{\alpha} = \mathbf{e} \cdot \mathbf{c}^{\alpha} = 1$ and thus $\mathbf{e} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{e}} = \frac{\partial (\mathbf{e} \cdot \mathbf{c}^{\alpha})}{\partial \mathbf{e}} = \mathbf{0}$, the identical argument being true for an arbitrary constant vector $\xi \mathbf{e}$ added to $\boldsymbol{\mu}$. Together, this shows that $\frac{1}{T} = \frac{1}{\theta}$, i.e. that the multiplier θ defining the local quasi-equilibrium has to be the same as the given parameter T .

- Inserting the definition of s into the definition (7.56), one has

$$f(\phi, \mathbf{c}, T) = \inf_{\mathbf{e}} \left\{ e - T \max_{(\mathbf{c}^{\alpha}, e^{\alpha})_{1 \leq \alpha \leq N} \in \mathcal{A}(\phi, \mathbf{c}, \mathbf{e})} \left\{ \sum_{\alpha} s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) h^{\alpha}(\phi) \right\} \right\}. \quad (7.58)$$

Since T is, as shown above, the same as the multiplier for the constraint $\sum_{\alpha=1}^N e^{\alpha} h^{\alpha}(\phi) = e$ in the local quasi-equilibrium conditions, the entropy can, given this value, equivalently be characterized through the free (with respect to the e^{α}) maximization problem

$$\begin{aligned} & \max_{(\mathbf{c}^{\alpha}, e^{\alpha})_{1 \leq \alpha \leq N} : (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N} \in \mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})} \left\{ \sum_{\alpha} s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) h^{\alpha}(\phi) - \frac{1}{T} \left(\sum_{\alpha=1}^N e^{\alpha} h^{\alpha}(\phi) - e \right) \right\} \\ &= \frac{e}{T} + \max_{(\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N} \in \mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})} \max_{(e^{\alpha})_{1 \leq \alpha \leq N}} \left\{ \sum_{\alpha} s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) h^{\alpha}(\phi) - \frac{1}{T} \sum_{\alpha=1}^N e^{\alpha} h^{\alpha}(\phi) \right\} \end{aligned}$$

where $\mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})$ is the admissible set from Equation (7.4) in the last section for defining f . Extracting the (fixed) factor $-\frac{1}{T}$ and using that the e^{α} are now independent, the last term can be rewritten as

$$\begin{aligned} & -\frac{1}{T} \min_{(\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N} \in \mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})} \min_{(e^{\alpha})_{1 \leq \alpha \leq N}} \left\{ \sum_{\alpha=1}^N \left(e^{\alpha} - T s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) \right) h^{\alpha}(\phi) \right\} \\ &= -\frac{1}{T} \min_{(\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N} \in \mathcal{A}_{\mathbf{c}}(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N \left(\min_{e^{\alpha}} \{ e^{\alpha} - T s^{\alpha}(\mathbf{c}^{\alpha}, e^{\alpha}) \} \right) h^{\alpha}(\phi) \right\}. \end{aligned}$$

The inner minimization problem is precisely the definition of the phase-specific free energy $f^\alpha(\mathbf{c}^\alpha, T)$ in terms of the phase-specific entropy s^α ,

$$f^\alpha(\mathbf{c}, T) = \min_{e^\alpha} \{e^\alpha - T s^\alpha(\mathbf{c}^\alpha, e^\alpha)\}, \quad (7.59)$$

which holds due to the “standard” thermodynamic relations valid for the bulk potentials, from which it follows that

$$\begin{aligned} & \max_{(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}: (\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha} s^\alpha(\mathbf{c}^\alpha, e^\alpha) h^\alpha(\phi) - \frac{1}{T} \left(\sum_{\alpha=1}^N e^\alpha h^\alpha(\phi) - e \right) \right\} \\ &= \frac{e}{T} - \frac{1}{T} \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\}. \end{aligned}$$

Inserting this expression into Equation (7.58), e in fact cancels out and one is left with

$$\begin{aligned} f(\phi, \mathbf{c}, T) &= \inf_e \left\{ e - T \left(\frac{e}{T} - \frac{1}{T} \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\} \right) \right\} \\ &= \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, T) h^\alpha(\phi) \right\}. \end{aligned}$$

The reverse conclusion is essentially just applying the same argument in the reverse order. Using Equation (7.59), one can rewrite f as

$$\begin{aligned} f(\phi, \mathbf{c}, T) &= \min_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N \min_{e^\alpha} \{e^\alpha - T s^\alpha(\mathbf{c}^\alpha, e^\alpha)\} h^\alpha(\phi) \right\} \\ &= \min_{e^\alpha} \left\{ \sum_{\alpha=1}^N \left\{ e^\alpha - T \max_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} s^\alpha(\mathbf{c}^\alpha, e^\alpha) \right\} h^\alpha(\phi) \right\}. \end{aligned}$$

The outer minimization in the e^α can be moved to the inner maximization problem by artificially “expanding” it through a double minimization

$$f(\phi, \mathbf{c}, T) = \min_e \min_{e^\alpha: \sum_{\alpha=1}^N e^\alpha h^\alpha(\phi) = e} \left\{ \sum_{\alpha=1}^N \left\{ e^\alpha - T \max_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} s^\alpha(\mathbf{c}^\alpha, e^\alpha) \right\} h^\alpha(\phi) \right\}$$

which is legitimate since, even though this adds an restriction on the sum of the previously completely free e^α , the value of sum can be chosen in any arbitrary manner to make the total expression as small as possible. Since the sum over the e^α appearing in the expression to be minimized is then obviously equal to e , one can move the \min_{e^α} -operation into the interior, leaving

$$f(\phi, \mathbf{c}, T) = \min_e \left\{ e - T \min_{e^\alpha: \sum_{\alpha=1}^N e^\alpha h^\alpha(\phi) = e} \max_{(\mathbf{c}^\alpha)_{1 \leq \alpha \leq N} \in \mathcal{A}_c(\phi, \mathbf{c})} \left\{ \sum_{\alpha=1}^N s^\alpha(\mathbf{c}^\alpha, e^\alpha) h^\alpha(\phi) \right\} \right\},$$

and the inner problem is just the definition of $s(\phi, \mathbf{c}, e)$ in Equation (7.52). \square

It turns out that the standard inversion formula for s in terms of f as

$$s(\phi, \mathbf{c}, e) = \inf_T \left\{ \frac{e}{T} - \frac{f(\phi, \mathbf{c}, T)}{T} \right\} \quad (7.60)$$

also continues to hold in this setting provided f is defined as in the previous section:

Lemma 8. Let $f(\phi, \mathbf{c}, T)$ be defined as in Equation (7.57) (resp. Equation (7.3)) and assume that the local quasi-equilibrium conditions (7.9) allow for a unique (up to the addition of an arbitrary constant vector $\xi \mathbf{e}$ to $\boldsymbol{\mu}$) solution which depends smoothly on \mathbf{c} and T and that the bulk free energy densities depends smoothly on \mathbf{c}^α and T . Then:

- Any minimizer T for the problem (7.60) is characterized by the averaging relation $e = \sum_{\alpha=1}^N e^\alpha(\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T)$ with respect to the external parameter e , where the $\mathbf{c}^\alpha(\phi, \mathbf{c}, T)$ are the ones defined in the optimality system for the definition of the f -function in Equation (7.3) and

$$e^\alpha(\mathbf{c}^\alpha, T) := -T^2 \frac{\partial}{\partial T} \left(\frac{f^\alpha(\mathbf{c}^\alpha, T)}{T} \right). \quad (7.61)$$

- The two definitions of $s(\phi, \mathbf{c}, e)$ through equations (7.52) and (7.60) coincide.

Proof. • Under the same simplifying smoothness assumptions on the solution of the local quasi-equilibrium condition from Equation (7.9), any minimizer T in the definition (7.60) has to satisfy

$$\begin{aligned} 0 &= -\frac{e}{T^2} + \sum_{\alpha=1}^N \frac{\partial}{\partial T} \left(\frac{f^\alpha(\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T)}{T} \right) h^\alpha(\phi) \\ &= -\frac{e}{T^2} + \frac{1}{T^2} \sum_{\alpha=1}^N \left(\frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha}(\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T) \frac{\partial \mathbf{c}^\alpha}{\partial T}(\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T) \right) h^\alpha(\phi) \\ &\quad - \frac{\partial}{\partial T} \left(\frac{f^\alpha}{T} \right) (\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T) h^\alpha(\phi), \end{aligned}$$

where the last derivative corresponds to the one of f^α with respect to its second argument only (i.e. while holding \mathbf{c}^α fixed).

Similarly to above, using $\frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha} = \boldsymbol{\mu} + \lambda^\alpha \mathbf{e}$ together with $\mathbf{e} \cdot \frac{\partial \mathbf{c}^\alpha}{\partial T} = 0$ and the independence of $\boldsymbol{\mu}$ on α shows that $\sum_{\alpha=1}^N \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha} \cdot \frac{\partial \mathbf{c}^\alpha}{\partial T} h^\alpha(\phi) = \boldsymbol{\mu} \cdot \frac{\partial \left(\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) \right)}{\partial T} = 0$ since the \mathbf{c}^α in the definition of $f(\phi, \mathbf{c}, T)$ satisfy the constraint $\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) = \mathbf{c}$. Using in addition the basic thermodynamic relation $-\frac{e^\alpha(\mathbf{c}, T)}{T} = \frac{\partial}{\partial T} \left(\frac{f^\alpha(\mathbf{c}^\alpha, T)}{T} \right)$ for the bulk free energy density, T thus has to be such that

$$-\frac{e}{T^2} + \sum_{\alpha=1}^N \frac{e^\alpha(\mathbf{c}^\alpha(\phi, \mathbf{c}, T), T)}{T^2} h^\alpha(\phi) = 0,$$

from which the claim follows.

- As was already seen in the proof of the previous lemma, the definition of $s(\phi, \mathbf{c}, e)$ in Equation (7.52) is such that the phase-specific temperatures T^α defined through $\frac{1}{T^\alpha} = \frac{\partial s^\alpha}{\partial e^\alpha}$ are equal to a common value θ and that the phase-specific concentrations are such that they coincide with the ones obtained from in the definition of $f(\phi, \mathbf{c}, \theta)$ in Equation (7.57) and thus correspond to the concentrations in the claim with T being replaced by θ . Furthermore, given any admissible phase-specific concentration \mathbf{c}^α and value of the multiplier T^α , it follows from “standard” thermodynamics that the values of the e^α can also be expressed as in equation (7.61)³⁸ in terms of \mathbf{c}^α and T^α . Evaluated for the particular choice $T^\alpha = \theta$ for all α , it follows from the sum-constraint on the e^α in the definition of (7.52) that these $e^\alpha(\mathbf{c}^\alpha, \theta)$ in fact satisfy $\sum_{\alpha=1}^N e^\alpha(\mathbf{c}^\alpha, \theta) = e$, and thus the condition in the first part of the lemma, from which it follows that $\theta = T$ under the assumption that there is a unique minimizer.

³⁸This is basically the phase-specific version of the conclusion obtained from Equation (7.56).

Conversely, any minimizer in Equation (7.60) does, by the use of $f(\phi, \mathbf{c}, T)$ through Equation (7.57) satisfy the equality of the chemical potentials. Furthermore, it follows from standard thermodynamics that defining $e^\alpha(\mathbf{c}^\alpha, T)$ as in Equation (7.61), that the $s^\alpha(\mathbf{c}^\alpha, e^\alpha(\mathbf{c}^\alpha, T))$ satisfy $\frac{1}{T} = \frac{\partial s}{\partial e^\alpha}(\mathbf{c}^\alpha, e^\alpha(\mathbf{c}^\alpha, T))$. Combined with the necessary condition for the minimizer in Equation (7.60), it follows that this choice of (\mathbf{c}^α, T) satisfies all the conditions required for the minimizer in Equation (7.52). \square

Remark 110. Note that the reasoning above relies mostly on a combination of two arguments, namely that standard thermodynamic relations holding at a phase-specific level and the optimality conditions for the phase-averaged quantities, ensuring that the contributions due to additional implicit dependencies drop out of the differential relations. It is likely that one can similarly extend most of the well-known “basic” thermodynamic relations to the phase-average setting in a very similar manner³⁹. \diamond

The Resulting Driving Force

The driving force for the phasefield equation can be derived in a manner very similar to the isothermal case in Section 7.1.3. The crucial point is again to keep in mind the dependence of the phase-specific concentrations and energies for the maximizers in Equation (7.52) on the phasefield ϕ . Assuming a smooth dependence and extending, if necessary, the optimality conditions in Equation (7.55) to those phases with $h^\alpha(\phi) = 0$, the differentiation of s with respect to ϕ^α leads to

$$\begin{aligned} \frac{\partial s(\phi, \mathbf{c}, T)}{\partial \phi^\alpha} &= \frac{\partial}{\partial \phi^\alpha} \left(\sum_\beta s^\beta(\mathbf{c}^\beta(\phi, \mathbf{c}, e), e^\beta(\phi, \mathbf{c}, e)) h^\beta(\phi) \right) \\ &= \sum_\beta s^\beta(\mathbf{c}^\beta(\phi, \mathbf{c}, e), e^\beta(\phi, \mathbf{c}, e)) \frac{\partial h^\beta}{\partial \phi^\alpha} + \sum_\beta \left(\frac{\partial s^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} + \frac{\partial s^\beta}{\partial e^\beta} \frac{\partial e^\beta}{\partial \phi^\alpha} \right) h^\beta(\phi). \end{aligned}$$

By construction, it holds that $\frac{\partial s^\beta}{\partial \mathbf{c}^\beta} = -\frac{\boldsymbol{\mu} + \lambda^\beta \mathbf{e}}{T}$ and $\frac{\partial s^\beta}{\partial e^\beta} = \frac{1}{T}$. Extracting the phase-independent quantities $\boldsymbol{\mu}$ and T out of the sum and making use of $\mathbf{e} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} = 0$ due to the sum-constraint $\mathbf{e} \cdot \mathbf{c}^\beta = 1$, it follows that

$$\sum_\beta \left(\frac{\partial s^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} + \frac{\partial s^\beta}{\partial e^\beta} \frac{\partial e^\beta}{\partial \phi^\alpha} \right) h^\beta(\phi) = -\frac{\boldsymbol{\mu}}{T} \cdot \sum_\beta \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi) + \frac{1}{T} \sum_\beta \frac{\partial e^\beta}{\partial \phi^\alpha} h^\beta(\phi).$$

The derivatives of the phase-specific quantities can be eliminated by an application of the product-rule since

$$\sum_\beta \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} h^\beta(\phi) = \frac{\partial \sum_\beta \mathbf{c}^\beta h^\beta(\phi)}{\partial \phi^\alpha} - \sum_\beta \mathbf{c}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} = \frac{\partial \mathbf{c}}{\partial \phi^\alpha} - \sum_\beta \mathbf{c}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} = -\sum_\beta \mathbf{c}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$$

and similarly $\sum_\beta \frac{\partial e^\beta}{\partial \phi^\alpha} h^\beta(\phi) = -\sum_\beta e^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$, leaving $\frac{\partial s(\phi, \mathbf{c}, T)}{\partial \phi^\alpha} = \sum_\beta \left(s^\beta + \frac{\boldsymbol{\mu}}{T} \mathbf{c}^\beta - \frac{1}{T} e^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha}$. This can further be simplified as $s^\beta - \frac{1}{T} e^\beta = -\frac{1}{T} (e^\beta - T s^\beta)$, whose value, as in Equation (7.59), coincides with $-\frac{1}{T} f^\beta(\mathbf{c}^\beta, T)$, now as a function of \mathbf{c}^β and T . This finally leads to

$$\frac{\partial s(\phi, \mathbf{c}, T)}{\partial \phi^\alpha} = \sum_\beta \left(-\frac{1}{T} f^\beta(\mathbf{c}^\beta, T) + \frac{\boldsymbol{\mu}}{T} \mathbf{c}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha} = -\frac{1}{T} \sum_\beta (f^\beta - \boldsymbol{\mu} \cdot \mathbf{c}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha}, \quad (7.62)$$

and thus again an expression based on the phase-specific grand chemical potential densities Ψ^β .

³⁹Recall also that it was already observed in Section 7.1.4 that $f(\phi, \mathbf{c}, T)$ and $\Omega(\phi, \boldsymbol{\mu}, T)$ are related through a standard Legendre transform.

The Solution Process

As already indicated in Remark 109 and as the discussion in the previous section shows, there in fact several different ways one can choose the “primary unknowns” and thus to solve the local quasi-equilibrium conditions in the non-isothermal case.

1. The first - and most obvious one - would be to apply a standard Newton-scheme directly to the original System (7.54), i.e. given the functions $s^\alpha(\mathbf{c}^\alpha, e^\alpha)$ and a set of values $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}^{(n)}$ for the phase-specific concentrations and energies as well as some value $(\boldsymbol{\eta}, \beta, (\hat{\lambda}^\alpha)_{1 \leq \alpha \leq N})^{(n)}$ for the associated multipliers, to solve for a set of corrections based on the linearization of Equation (7.54). By linearity in the multipliers $(\boldsymbol{\eta}, \beta, (\hat{\lambda}^\alpha)_{1 \leq \alpha \leq N})$, this can, as in Section 7.1.5, in effect be reduced to an iteration based on the $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}^{(n)}$

alone, with, again by linearity, the residuals in the sum-constraints $\sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\boldsymbol{\phi}) = \mathbf{c}$ and $\sum_{\alpha=1}^N e^\alpha h^\alpha(\boldsymbol{\phi}) = e$ vanishing after the first iteration. This has the advantage of on the one hand requiring only direct evaluations of the derivatives of the (assumed to be given) $s^\alpha(\mathbf{c}^\alpha, e^\alpha)$ and on the other hand ensuring “conservation” of the phase-averaged conserved (\mathbf{c}, e) variables after a single Newton-step, regardless of the accuracy of the estimates for the multipliers.

The major disadvantage is that this on the one hand requires working with the - less commonly used - functions $s^\alpha(\mathbf{c}^\alpha, e^\alpha)$ and on the other hand leads to a system with total of $(K + 1)N$ “primary” unknowns $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ in addition to the $K + N + 1$ multipliers $(\boldsymbol{\eta}, \beta, (\hat{\lambda}^\alpha))$ (resp. KN primary unknowns $(\tilde{\mathbf{c}}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ and K unknown multipliers $(\tilde{\boldsymbol{\mu}}^\alpha, \beta)$ after a reduction to $K - 1$ independent concentration values as in Section 7.1.5).

2. A quite natural alternative is suggested by the discussion in Section 7.1.5 in the isothermal case. One can directly make use of the equilibrium condition $T^\alpha = T$ and then to use the equivalent System (7.55), but instead based on the $(\mathbf{c}^\alpha, T^\alpha)_{1 \leq \alpha \leq N}$ as primary unknowns. One advantage of this approach is that, since the local quasi-equilibrium conditions from Equation (7.9) in the isothermal case are implicitly also based upon the assumption that each T^α corresponds to the (given) value T , the system in Equation (7.55) can basically be considered to be the same as in the previous case. The only difference is that now the temperature T forms part of the unknowns and has to be fixed such that the averaging condition $\sum_{\alpha=1}^N e^\alpha(\mathbf{c}^\alpha, T)h^\alpha(\boldsymbol{\phi}) = e$ - now with the e^α as explicit functions of \mathbf{c}^α and T - holds.

A second one is that, as the trivial constraint $T^\alpha = T$ can directly be integrated into the system. As the $(e^\alpha)_{1 \leq \alpha \leq N}$ are here considered to be known for given values of (\mathbf{c}^α, T) , one can reduce the number of primary unknowns to the $KN + 1$ values $((\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, T)$ and $K + N$ additional secondary unknowns through the multipliers $(\boldsymbol{\mu}, (\lambda^\alpha)_{1 \leq \alpha \leq N})$. Using a reduced formulation, further eliminates unknowns, leaving only $(K - 1)N + 1$ primary unknowns $((\tilde{\mathbf{c}}^\alpha)_{1 \leq \alpha \leq N}, T)$ and the $(K - 1)$ values of the multiplier $\tilde{\boldsymbol{\mu}}$.

One disadvantage - if based on storing (\mathbf{c}, T) as the primary unknowns instead of (\mathbf{c}, e) - is that one now needs to ensure that the (local) quasi-equilibrium systems are solved accurately in order to ensure conservation of energy. This is due to the fact that the function $e^\alpha(\mathbf{c}^\alpha, T)$ are now in general nonlinear in both \mathbf{c}^α and T . The same then holds for the the sum-constraint $\sum_{\alpha=1}^N e^\alpha(\mathbf{c}^\alpha, T) = e$, which will therefore not be solved exactly with a single correction step⁴⁰.

⁴⁰Note that this can arise even for quite simple dependencies of the energies on the concentration and temperature including e.g. a product of a linearly (in \mathbf{c}^α) interpolated specific heat capacity $(c_v)^\alpha(\mathbf{c}^\alpha)$ with the temperature. Even both factor are linear, the product obviously is not.

3. A third approach - which is closely related to the approach suggested in [19] - is to use a fully “dual” description, i.e. by also replacing the \mathbf{c}^α as functions of $(\boldsymbol{\mu}, (\lambda^\alpha)_{1 \leq \alpha \leq N}, T)$ and (thus also, through the \mathbf{c}^α , the e^α as functions of the same parameters) and then applying a Newton-scheme on this set of equations. A clear advantage is that this leads to an even further reduction of the primary unknowns to the $K + N + 1$ values (resp. the K values $(\tilde{\boldsymbol{\mu}}, T)$ when using a reduced version) as the $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ are now taken as direct functions of the $\boldsymbol{\mu}^\alpha = \boldsymbol{\mu} + \lambda^\alpha \mathbf{e}$ and T (resp. $\tilde{\boldsymbol{\mu}}$ and T).

Which approach is more appropriate again depends quite heavily on which thermodynamic potential(s) can be considered as “given”. Even though the increasing reduction in the number of unknowns a priori seems highly favorable, it is subject to the same pitfalls as the choice of working with $\tilde{\boldsymbol{\mu}}$ as the primary unknown already discussed in Remark 104. In particular, if based on $s^\alpha(\mathbf{c}^\alpha, e^\alpha)$ as the basic thermodynamic potential, using $((\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, T)$ instead of $(\mathbf{c}^\alpha, e^\alpha)_{1 \leq \alpha \leq N}$ as the primary unknowns is in general highly inconvenient as the solution of the subsystems $\frac{\partial s^\alpha}{\partial e^\alpha}(\mathbf{c}^\alpha, e^\alpha) = \frac{1}{T}$ for determining the (then dependent) e^α in terms of the (assumed given) \mathbf{c}^α and T can itself require a number of Newton-steps. Even though this is scalar equation in a single unknown, this is still likely an effort ill spent. One can of course then proceed similarly to Remark 104 and allow for controlled inaccuracies, this leading, similarly to the isothermal case based on $\tilde{\boldsymbol{\mu}}$ as the primary unknown, essentially back to first approach. The situation is even worse when taking $(\tilde{\boldsymbol{\mu}}^\alpha, T)$ as the primary unknowns, as this then requires, for each evaluation of $\frac{\partial^2 s^\alpha}{\partial (\mathbf{c}^\alpha)^2}$, $\frac{\partial^2 s^\alpha}{\partial \mathbf{c}^\alpha \partial e^\alpha}$ and $\frac{\partial^2 s^\alpha}{\partial (e^\alpha)^2}$ solving the systems

$$\frac{\partial s^\alpha}{\partial e^\alpha}(\tilde{\mathbf{c}}^\alpha, e^\alpha) \stackrel{!}{=} \frac{1}{T} \quad \text{and} \quad \frac{\partial s^\alpha}{\partial \tilde{\mathbf{c}}^\alpha}(\tilde{\mathbf{c}}^\alpha, e^\alpha) \stackrel{!}{=} \frac{\tilde{\boldsymbol{\mu}}^\alpha}{T}$$

to a sufficiently high accuracy.

If based on $f^\alpha(\mathbf{c}^\alpha, T)$ as the basic thermodynamic potential, as is currently the standard in the **Pace3D**-framework, the use of $((\mathbf{c}^\alpha)_{1 \leq \alpha \leq N}, T)$ is in general the most convenient one, and will therefore be the one considered here.

Given that the fundamental conservative evolution equations are based on the concentration \mathbf{c} and the energy e , the requirement is therefore to find phase-specific concentrations \mathbf{c}^α and energies e^α compatibly with the optimality conditions for Equation (7.52) given the values of the average concentration \mathbf{c} and the average energy e . Since it is the f^α as functions of \mathbf{c}^α and T^α which are, for practical reasons, chosen as the fundamental potentials, the natural parameterization for doing so is also in terms of the concentration values \mathbf{c}^α and the - known to have to be a common value for all α - temperature T .

Given some initial guess $\mathbf{c}^{(0)}$ and $T^{(0)}$, it is again natural to use a Newton-scheme for the solution of the optimality system in Equation (7.55).

Remark 111. As already indicated, a parameterization in terms of T instead of e in principle has the same drawbacks already discussed in Remarks 104 and 105. Being based on f^α formulated directly in terms of T , the constructions are usually such that the conversion of a given concentration and temperature to the corresponding energy is relatively simple, and thus less problematic. It will therefore be assumed in the following that obtaining the phase-specific energies e^α as a function of (\mathbf{c}^α, T) is computationally relatively cheap.

If such is not the case, one may have to apply similar considerations as in Remark 104. \diamond

Given some initial guess such e.g. $(\mathbf{c}^\alpha)^{(0)} = \mathbf{c}$, $(e^\alpha)^{(0)} = e^\alpha((\mathbf{c}^\alpha)^{(0)}, T^{(0)})$, the basic Newton-

step consists in solving

$$\begin{cases} \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) \delta \mathbf{c}^\alpha + \frac{\partial \bar{\mu}^\alpha}{\partial T} \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) \delta T - \bar{\mu}^{(n+1)} - (\lambda^\alpha)^{(n+1)} \mathbf{e} \right) \stackrel{!}{=} -\bar{\mu}^\alpha \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) \\ \sum_{\alpha=1}^N \delta \bar{\mathbf{c}}^\alpha h^\alpha \stackrel{!}{=} \mathbf{c} - \sum_{\alpha=1}^N (\mathbf{c}^\alpha)^{(n)} h^\alpha =: \mathbf{r}_{\bar{\mathbf{c}}} \\ \sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial \bar{c}^\alpha} \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) \delta \mathbf{c}^\alpha + \frac{\partial e^\alpha}{\partial T} \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) \delta T \right) h^\alpha \stackrel{!}{=} e - \sum_{\alpha=1}^N e^\alpha \left((\mathbf{c}^\alpha)^{(n)}, T^{(n)} \right) =: r_e \end{cases}$$

Due to the same considerations as in Section 7.1.5, this system is again most easily solved by reverting to a reduced form, either by a priori using a reduced formulation or using an algebraic “a posteriori” reduction. Dropping the arguments for notational simplicity, this results in having to solve the simpler system⁴¹

$$\begin{cases} \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \delta \bar{\mathbf{c}}^\alpha + \frac{\partial \bar{\mu}^\alpha}{\partial T} \delta T - \tilde{\mu}^{(n+1)} \right) \stackrel{!}{=} -(\tilde{\mu}^\alpha)^{(n)} \\ \sum_{\alpha=1}^N \delta \bar{\mathbf{c}}^\alpha h^\alpha \stackrel{!}{=} \mathbf{c} - \sum_{\alpha=1}^N (\mathbf{c}^\alpha)^{(n)} h^\alpha = \mathbf{r}_{\bar{\mathbf{c}}} \\ \sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial \bar{c}^\alpha} \delta \bar{\mathbf{c}}^\alpha + \frac{\partial e^\alpha}{\partial T} \delta T \right) h^\alpha \stackrel{!}{=} e - \sum_{\alpha=1}^N (e^\alpha)^{(n)} = r_e. \end{cases} \quad (7.63)$$

Since this system is still in a block-diagonal form with respect to the $\delta \bar{\mathbf{c}}^\alpha$, it is again convenient to perform a block-elimination in this equation, leading, completely analogous to the discussion in Section 7.1.5, to the Schur-complement system

$$\begin{cases} \left(\sum_{\alpha=1}^N \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} \left(\tilde{\mu}^{(n+1)} - (\tilde{\mu}^\alpha)^{(n)} - \frac{\partial \bar{\mu}^\alpha}{\partial T} \delta T \right) \right) h^\alpha \stackrel{!}{=} r_{\bar{\mathbf{c}}} \\ \left(\sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial \bar{c}^\alpha} \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} \left(\tilde{\mu}^{(n+1)} - (\tilde{\mu}^\alpha)^{(n)} - \frac{\partial \bar{\mu}^\alpha}{\partial T} \delta T \right) + \frac{\partial e^\alpha}{\partial T} \delta T \right) \right) h^\alpha \stackrel{!}{=} r_e, \end{cases} \quad (7.64)$$

consisting of a total of K unknowns, $K - 1$ ones in terms of $\tilde{\mu}$ and one in terms of δT .

K is, due to computational requirements, usually quite small (typically below 4 or 5) due to the fact that the concentration evolution - unlike an obstacle-potential based phasefield equation - has to be calculated everywhere. A very natural choice for determining the values of $\delta \bar{\mu}$ and δT in Equation (7.64) is therefore simply to use any direct solver, since this is a non-singular system of typically very small size. Given the values of $\tilde{\mu}$ and δT , it is then again an easy matter of updating the remaining phase-specific quantities $\delta \bar{\mathbf{c}}^\alpha$ based on Equation (7.63), in which all phases are then uncoupled.

Once the iterations converge, one can then, if required recover the values of $\bar{\mu}$ and the λ^α as in Subsection 7.1.5, or, if not so, simply recover the concentration values c_K^α based on the sum-constraint $\sum_{i=1}^K c_i^\alpha = 1$.

Remark 112. In terms of code-reusability, it may nevertheless have some interest to, in the same spirit as in Subsection 7.1.5, perform another block-reduction on Equation (7.64). First moving all known quantities to the right-hand side, this system can be rewritten as

$$\begin{cases} \left(\mathbf{S}_{\bar{\mathbf{c}}} \tilde{\mu}^{(n+1)} - \left(\sum_{\alpha=1}^N \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} \frac{\partial \bar{\mu}^\alpha}{\partial T} h^\alpha \right) \delta T \right) = \mathbf{r}_{\bar{\mathbf{c}}} + \sum_{\alpha=1}^N \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} (\tilde{\mu}^\alpha)^{(n)} h^\alpha =: \hat{\mathbf{r}}_{\bar{\mathbf{c}}}, \\ \left(\left(\sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial \bar{c}^\alpha} \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} h^\alpha \right) \tilde{\mu}^{(n+1)} + S_T \delta T \right) = r_e + \sum_{\alpha=1}^n \frac{\partial e^\alpha}{\partial \bar{c}^\alpha} \left(\frac{\partial \bar{\mu}^\alpha}{\partial \bar{c}^\alpha} \right)^{-1} (\tilde{\mu}^\alpha)^{(n)} h^\alpha =: \hat{r}_e, \end{cases} \quad (7.65)$$

⁴¹Introducing the reduction into the term $\frac{\partial e^\alpha}{\partial \bar{c}^\alpha}$ is equally simple as in Subsection 7.1.5 using $c_K^\alpha = 1 - \sum_{i=1}^{K-1} c_i^\alpha$.

where $\mathbf{S}_{\tilde{\mathbf{c}}}$ is the same Schur-complement matrix as in Equation (7.37) in Section 7.1.5 and S_T is the analogous “matrix” (this is simply a scalar value) to $\mathbf{S}_{\tilde{\mathbf{c}}}$ for the temperature increment given by

$$S_T = \sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial T} - \frac{\partial e^\alpha}{\partial \mathbf{c}^\alpha} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial T} \right) h^\alpha. \quad (7.66)$$

The two basic choices for this are to either make use of the first equation in Equation (7.64) for eliminating $\tilde{\boldsymbol{\mu}}^{(n+1)}$ as a function of δT ,

$$\tilde{\boldsymbol{\mu}}^{(n+1)}(\delta T) = \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \hat{\mathbf{r}}_{\tilde{\mathbf{c}}} + \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \sum_{\alpha=1}^N \left(\left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial T} h^\alpha \right) \delta T$$

or the elimination of δT as a function of $\tilde{\boldsymbol{\mu}}^{(n+1)}$ using the second equation, leading to

$$\delta T(\tilde{\boldsymbol{\mu}}^{(n+1)}) = \frac{1}{S_T} \hat{r}_e - \frac{1}{S_T} \left(\sum_{\alpha=1}^n \frac{\partial e^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} h^\alpha \right) \tilde{\boldsymbol{\mu}}^{(n+1)}.$$

Inserting these two expressions into the respective other equation in the System (7.65), the first choice leads a scalar equation in δT with a modified Schur-complement

$$S_T + \left(\sum_{\alpha=1}^n \frac{\partial e^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} h^\alpha \right) \mathbf{S}_{\tilde{\mathbf{c}}}^{-1} \sum_{\alpha=1}^N \left(\left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial T} h^\alpha \right)$$

for the temperature increment, whereas the second choice leads to a modified equation for $\tilde{\boldsymbol{\mu}}^{(n+1)}$ on the modified Schur-complement

$$\mathbf{S}_{\tilde{\mathbf{c}}} + \frac{1}{S_T} \left(\sum_{\alpha=1}^N \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial T} h^\alpha \right) \left(\sum_{\alpha=1}^n \left(\frac{\partial e^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)^{-1} h^\alpha \right)$$

for the reduced chemical potential.

If the code-functionality is primarily designed with the isothermal case in mind, the former choice is likely the more convenient one, since this makes it possible to perform the solution of the System (7.63) essentially without modification to the “purely chemical” part, as one can reuse this functionality in an exterior manner for determining δT . Once δT is known, the remaining system in $\tilde{\boldsymbol{\mu}}^{(n+1)}$ and the $\delta \tilde{\mathbf{c}}^\alpha$ then reduces to the one in Section 7.1.5 with slightly modified right-hand sides.

A similar construction with a larger number of additional unknowns will be discussed in more detail in Section 7.2.5. \diamond

Remark 113. Note that it is also possible to obtain slightly different expressions making use of some thermodynamic relations. For example, based on $e^\alpha(\mathbf{c}^\alpha, T) = -T^2 \frac{\partial}{\partial T} \left(\frac{f^\alpha(\mathbf{c}^\alpha, T)}{T} \right)$, one has

$$\begin{aligned} \frac{\partial e^\alpha}{\partial \mathbf{c}^\alpha} &= -T^2 \frac{\partial^2}{\partial \mathbf{c}^\alpha \partial T} \left(\frac{f^\alpha(\mathbf{c}^\alpha, T)}{T} \right) = -T^2 \frac{\partial}{\partial T} \left(\frac{1}{T} \frac{\partial f^\alpha}{\partial \mathbf{c}^\alpha} \right) = -T^2 \frac{\partial}{\partial T} \left(\frac{\boldsymbol{\mu}^\alpha(\mathbf{c}^\alpha, T)}{T} \right) \\ &= -T^2 \left(-\frac{\boldsymbol{\mu}^\alpha(\mathbf{c}^\alpha, T)}{T^2} + \frac{1}{T} \frac{\partial \boldsymbol{\mu}^\alpha}{\partial T} \right) = \boldsymbol{\mu}^\alpha(\mathbf{c}^\alpha, T) - T \frac{\partial \boldsymbol{\mu}^\alpha}{\partial T}. \end{aligned}$$

In terms of “readability”, an obvious definition to make use of is $c_v^\alpha = \frac{\partial e^\alpha}{\partial T}$. \diamond

7.2 Quantitative Models for Solid-Solid Transformations

Similarly to the solid-liquid transformations considered in the previous section, the phasefield method has also gained considerable interest for simulating the evolution of microstructures consisting of different solid phases. Such solid-solid state transformations are, in addition to the surface energy of the various phases, also heavily influenced by stored mechanical energy contributions.

In some cases, for example during recrystallization processes, this additional stored energy distribution can be reasonably approximated as being a function of the phase only, thus allowing for a simpler model as in Section 6 with the phasefields being the only unknowns (see e.g. [78]). In other cases, the phase-transformation process itself heavily influences the mechanical state of the system, requiring a more complex treatment in terms of a coupled phasefield-mechanical model, where now there is (at least) one additional set of unknowns in terms of the displacement field \mathbf{u} .

Phasefield models for solid-solid transformations are usually again based upon an appropriately chosen energy functional in combination with a gradient-flow postulate on the dynamics of the phasefield variables. In the small deformation setting, the phasefield functional is therefore usually chosen (see e.g. [69] and [3] as some of the earlier works coupling phasefield methods with elasticity) as consisting of the standard phasefield contributions through a and w supplemented by the total elastic energy

$$\mathcal{F}_\epsilon(\phi, \mathbf{u}) = \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) \, d\mathbf{x} + \mathcal{F}_{el}(\phi, \mathbf{u}) \quad (7.67)$$

where, assuming for example a mixture of displacement and traction boundary conditions on the Dirichlet-part Γ_D and the Neumann-part Γ_N , $\mathcal{F}_{el}(\phi, \mathbf{u})$ is given by⁴²

$$\mathcal{F}_{el}(\phi, \mathbf{u}) = \int_{\Omega} f_{el}(\phi, \boldsymbol{\epsilon}(\mathbf{u})) \, d\mathbf{x} - \rho \mathbf{f} \cdot \mathbf{u} \, d\mathbf{x} - \int_{\Gamma_N} \mathbf{g} \cdot \mathbf{u} \, ds \quad (7.68)$$

with the (here volumetric) strain energy density $f_{el}(\phi, \boldsymbol{\epsilon}(\mathbf{u}))$ depending on \mathbf{u} only through the symmetric displacement gradient $\boldsymbol{\epsilon}(\mathbf{u}) = \nabla_S \mathbf{u}$.

In contrast to the solidification models considered in Section 7.1, the evolution of the mechanical state is not driven by a (typically comparatively slow) diffusion process, but by a wave-propagation process which tends to quickly equilibrate itself until a mechanical equilibrium is reached. As this equilibration often takes place at a much shorter timescale than that of the evolution of the phase boundaries⁴³, it is common to employ a quasi-static approximation for the mechanical fields.

Remark 114. While it is well-known that the propagation of a perfectly elastic wave conserves energy, there are additional dissipative effects such as the radiation of wave energy into the environment and damping of the waves within the material. These will in reality ultimately lead to the establishment of a mechanically equilibrated state. The validity of the quasi-static approximation therefore depends upon how fast this equilibration takes place as compared to the other processes involved. \diamond

⁴²Note that the contribution due to \mathbf{f} and the boundary contribution is often left out, even when considering loads through body forces and/or external stresses. While these do not change the resulting phasefield evolution equation if one does a purely formal differentiation (i.e. “forgetting” the coupling of \mathbf{u} to ϕ), it is crucial for the actual validity of this calculation.

⁴³This need not always be the case though as for example martensitic phase transformations occur at a very fast rate.

Under this quasi-static assumption, \mathbf{u} is subject to the mechanical equilibrium condition

$$\begin{cases} -\operatorname{div}(\boldsymbol{\sigma}(\boldsymbol{\phi}, \boldsymbol{\epsilon}(\mathbf{u}))) = \mathbf{f} & \text{in } \Omega, \\ \mathbf{u} = \mathbf{U} & \text{on } \Gamma_D, \\ \boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{g} & \text{on } \Gamma_N \end{cases} \quad (7.69)$$

where $\boldsymbol{\sigma}$ denotes the Cauchy stress tensor. Note that in order to be consistent with a variational principle based on the functional $\mathcal{F}_\epsilon(\boldsymbol{\phi}, \mathbf{u})$ from Equation (7.67) and with Equation (7.68), $\boldsymbol{\sigma}(\boldsymbol{\phi}, \boldsymbol{\epsilon}(\mathbf{u}))$ in addition has to satisfy (see the discussion below)

$$\boldsymbol{\sigma}(\boldsymbol{\phi}, \boldsymbol{\epsilon}(\mathbf{u})) = \frac{\partial f_{el}(\boldsymbol{\phi}, \boldsymbol{\epsilon})}{\partial \boldsymbol{\epsilon}}(\boldsymbol{\phi}, \boldsymbol{\epsilon}(\mathbf{u})). \quad (7.70)$$

As the external loads through \mathbf{f} and \mathbf{g} are to be considered as given data, it is clear that the only point in which the various phasefield models differ lies in the definition of this strain energy density f_{el} . Alternatively, and which amounts up to an integration constant to the same provided⁴⁴ $\mathcal{C}(\boldsymbol{\phi}, \boldsymbol{\epsilon}) := \frac{\partial \boldsymbol{\sigma}}{\partial \boldsymbol{\epsilon}}$ obeys the usual symmetry conditions compatible with $\mathcal{C}(\boldsymbol{\phi}, \boldsymbol{\epsilon}) = \frac{\partial^2 f_{el}}{\partial \boldsymbol{\epsilon}^2}$, one may also define a stress-strain relationship $\boldsymbol{\sigma}(\boldsymbol{\phi}, \boldsymbol{\epsilon})$.

Following the former approach and allowing for **eigenstrains (prestrains)** $\tilde{\boldsymbol{\epsilon}}^\alpha$, this strain energy density within an α -bulk region is given by (see e.g. [69])

$$f_{el}^\alpha(\boldsymbol{\epsilon}^\alpha) = \frac{1}{2} \boldsymbol{\epsilon}_{el}^\alpha : \mathcal{C}^\alpha : \boldsymbol{\epsilon}_{el}^\alpha = \frac{1}{2} (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) : \mathcal{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha),$$

where the elastic strain $\boldsymbol{\epsilon}_{el}^\alpha$ is defined as the difference of the total and prestrains, $\boldsymbol{\epsilon}_{el}^\alpha = \boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha$. As for the solidification problems considered in the previous section, the primary difficulty for designing an accurate phasefield model thus lies in properly extending this definition to within the interface regions. Based on the bulk-expression, the natural approach within the phasefield context is again the use of a weighted interpolation

$$f_{el}(\boldsymbol{\phi}, (\boldsymbol{\epsilon}^\alpha)_{1 \leq \alpha \leq N}) = \sum_{\alpha=1}^N f_{el}^\alpha(\boldsymbol{\epsilon}^\alpha) h^\alpha(\boldsymbol{\phi}) = \sum_{\alpha=1}^N \left(\frac{1}{2} (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) : \mathcal{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) \right) h^\alpha(\boldsymbol{\phi}) \quad (7.71)$$

of the elastic free energy contributions of the individual phases. This is still an incomplete description though as it remains to specify the dependence of the phase-specific strains $\boldsymbol{\epsilon}^\alpha$ on the total strain $\boldsymbol{\epsilon}$.

One natural restriction is that one would like for the relation $\boldsymbol{\sigma}^\alpha = \mathcal{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ to hold. A second one - which is typically imposed in accordance with the standard approach of representing the total quantities as an interpolation of the various phases - is that the total strain and the total stresses be related to the phase-specific ones through

$$\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^\alpha h^\alpha(\boldsymbol{\phi}) \quad \text{and} \quad \boldsymbol{\sigma} = \sum_{\alpha=1}^N \boldsymbol{\sigma}^\alpha h^\alpha(\boldsymbol{\phi}). \quad (7.72)$$

This still leaves a large degree of flexibility in the construction of an effective material behavior though.

⁴⁴Even though, at least in the linearly elastic case, it is natural here to assume that \mathcal{C} is independent of $\boldsymbol{\epsilon}$ (i.e. that the diffuse interface stress-strain relationship is linear in $\boldsymbol{\epsilon}$), this is not strictly necessary. While a nonlinear stress-strain relationship a priori would then violate the linearity of the original problem, this could still be of interest if this “violation” is carefully designed. One example of a similar situation is the use of nonlinear advection schemes for an a priori linear advection operator.

In terms of its contribution to the phasefield Equation (6.73), the minimization of \mathcal{F}_ϵ in Equation (7.67) subject to the quasi-static equilibrium conditions (7.69) a priori leads to a constrained optimization problem in terms of ϕ and \mathbf{u} . Equation (7.69) is, under mild hypothesis, uniquely (or uniquely up to some irrelevant kernel) solvable for \mathbf{u} as a function of ϕ . One can thus again consider a reduced formulation consisting in the minimization of $\mathcal{F}_\epsilon(\phi, \mathbf{u}(\phi))$ in terms of which the FONC for any minimizer becomes

$$\left\langle \frac{\partial \mathcal{F}_\epsilon}{\partial \phi}, \delta \phi \right\rangle + \left\langle \frac{\partial \mathcal{F}_\epsilon}{\partial \mathbf{u}}, \delta \mathbf{u}(\delta \phi) \right\rangle = \langle \boldsymbol{\mu} + \Lambda, \delta \phi \rangle$$

with $\boldsymbol{\mu}$ and Λ as in Section 6.1.

As \mathbf{u} only enters in $\mathcal{F}_{el}(\phi, \mathbf{u})$ and the system (7.69) is, under the variational consistency condition (7.70), precisely the condition $\langle \frac{\partial \mathcal{F}_{el}(\phi, \mathbf{u})}{\partial \mathbf{u}}, \delta \mathbf{u} \rangle = 0$ for all admissible variations $\delta \mathbf{u}$. The second contribution therefore actually vanishes and one only has to consider the additional “driving force” through the direct contribution $\frac{\partial \mathcal{F}_{el}}{\partial \phi}$ to the phasefield equation. Assuming for simplicity that neither the body force $\rho \mathbf{f}$ nor the imposed boundary conditions depend on ϕ , the only explicit dependence on ϕ is therefore through the strain energy density $f_{el}(\phi, \boldsymbol{\epsilon}(\mathbf{u}))$ in Equation (7.71). In the simpler models (this will be called the “traditional” phasefield models below), the dependence of f_{el} on ϕ is solely in terms of the (local) values of ϕ itself and not its gradients. From Equation (7.71) and the chain-rule, this leads to an additional contribution of the form

$$\frac{\partial f_{el}}{\partial \phi^\alpha} = \sum_{\beta=1}^N f^\beta(\boldsymbol{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha}(\phi) + \sum_{\beta=1}^N h^\beta(\phi) \frac{\partial f^\beta}{\partial \boldsymbol{\epsilon}^\beta} : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \phi^\alpha}.$$

For the more complex models discussed starting from Section 7.2.2, the calculation of f_{el} relies heavily on one or several normal vectors $\mathbf{n}^{\alpha\beta}(\phi, \nabla \phi)$ between the various phases, defined in terms of (possibly) the local values ϕ and their gradients $\nabla \phi$. It is therefore more convenient to write f_{el} and $\boldsymbol{\epsilon}^\alpha$ as functions of $\epsilon^{45}(\boldsymbol{\epsilon}, \phi, \mathbf{n}^{\alpha\beta}(\phi, \nabla \phi))_{1 \leq \alpha \neq \beta \leq N}$, from which one obtains the additional contributions

$$\sum_{\beta=1}^N \sum_{\delta=1}^N \left(h^\beta(\phi) \frac{\partial f^\beta}{\partial \boldsymbol{\epsilon}^\beta} : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \mathbf{n}^{\beta\delta}} \cdot \frac{\partial \mathbf{n}^{\beta\delta}}{\partial \phi^\alpha} - \nabla \cdot \left(h^\beta(\phi) \frac{\partial f^\beta}{\partial \boldsymbol{\epsilon}^\beta} : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \mathbf{n}^{\beta\delta}} \cdot \frac{\partial \mathbf{n}^{\beta\delta}}{\partial \nabla \phi^\alpha} \right) \right).$$

Remark 115. Note the a priori somewhat curious situation that, even though the interpretation of the $\boldsymbol{\epsilon}^\alpha$ is as the phase-inherent strains for the given phase α and thus in particular not a “phase-averaged” quantity, these will, except for the simplest models, still depend implicitly upon (at least) the ϕ -values. In contrast, the total strain $\boldsymbol{\epsilon}$ - even though actually interpreted as the phase-averaged quantity - is to be considered as independent of the phasefield (and the other parameters).

This situation is completely analogous to Section 7.1, where the average concentration \mathbf{c} was considered as phase-independent, whereas the phase-specific concentration values depended upon ϕ through their specification in terms of the condition of equal (reduced) chemical potentials. This is clearly an inherent feature of the common approach underlying these phasefield models, which, for obvious efficiency reasons, consists in expressing the phasefield functional and thus the corresponding equilibrium conditions through some variational principle in terms of the phase-averaged quantities as the primary unknown. The phase-specific quantities are only introduced as auxiliary “secondary” quantities through which the bulk free energy potential in the interface region is then (locally) linked to the ones of the individual phases. This is also why there is - except potentially for the bulk-regions - no reason to expect the phase-inherent quantities to

⁴⁵Since both the definition of suitable normal vectors in the multiphase case and the definition of normals at the transition to bulk phases with $\nabla \phi = \mathbf{0}$ are by themselves somewhat tricky issue, the contributions due to the dependences of the normal vectors will later on mostly be expressed only in terms of $\frac{\partial \boldsymbol{\epsilon}^\alpha}{\partial \mathbf{n}^{\beta\delta}}$, leaving the specifics of the definition of those aside.

obey the same elegant relations holding for the phase-averaged ones. For example, the ϵ^α defined by the various models will in general not derive from a global displacement field \mathbf{u}^α and neither will the corresponding phase-specific equilibrium condition $-\text{div}(\boldsymbol{\sigma}^\alpha) = \rho^\alpha f^\alpha$ hold throughout the domain (but it will, of course, be valid in the respective bulk-region). \diamond

Remark 116. It is of course also possible to consider situations in which the density and/or body force as well as the boundary conditions in addition depend on ϕ . This is easily dealt with in the former case as, typically being defined in a simple local relation in terms of ϕ , the contribution to the phasefield equation follows from a straightforward differentiation of this dependence. In contrast, the latter case can, depending upon the complexity of the boundary condition, become significantly more technical to deal with.

If \mathbf{g} for example only depends on ϕ in a simple local manner such as e.g. when imposing \mathbf{g} as an interpolation of normal stresses assigned to the individual phases, this only affects the phasefield equation in the sense that the natural “isolating” boundary condition needs to be replaced by one including $\frac{\partial \mathbf{g}}{\partial \phi^\alpha}$ ⁴⁶. In other cases - quite intuitive from a purely mechanical point of view - such as e.g. for a “constant-force” boundary condition, where a given force is distributed over the boundary based upon some interpolation procedure in terms of ϕ , one has to deal with a non-local dependence of \mathbf{g} on ϕ on the Neumann-part of the boundary. This leads to a correspondingly more complex (non-local) contribution to the phasefield boundary condition. \diamond

7.2.1 “Traditional” Phasefield Models

Before considering a more recent modeling approach in a little more detail, this section will recall a few key points concerning three more classical modeling approaches, namely the Voigt-Taylor-model, the Reuss-Sachs-model and the model by Khachaturian. As the differences between these models have been amply analyzed in the literature (see e.g. [3], [23], [51] and [64]), the primary purpose of this section is to provide a brief summary of this analysis (and how it relates to the one in the previous section) as a background for the following discussion.

The Purely Elastic Case

Two popular early approaches are to either assume that the strains or the stresses of all phases are equal. In the first approach, the so-called **Voigt-Taylor model**, the assumption $\epsilon^\alpha = \epsilon \forall \alpha$, the first relation in Equation (7.72) is trivially satisfied. Based on imposing $\boldsymbol{\sigma}^\alpha = \mathbf{C}^\alpha : \epsilon^\alpha$ and the second relation in (7.72), this leads to the total stress

$$\boldsymbol{\sigma}_{VT}(\phi, \epsilon) = \sum_{\alpha=1}^N \boldsymbol{\sigma}^\alpha(\epsilon) h^\alpha(\phi) = \left(\sum_{\alpha=1}^N \mathbf{C}^\alpha h^\alpha(\phi) \right) : \epsilon, \quad (7.73)$$

i.e. $\boldsymbol{\sigma}$ is determined in terms of the total strain through an effective stiffness tensor given by the (weighted) arithmetic interpolation of the phase-specific ones, $\mathbf{C}_{VT}(\phi) = \sum_{\alpha=1}^N \mathbf{C}^\alpha h^\alpha(\phi)$.

In addition, the total local free energy contribution is of the form

$$f_{el,VT}(\phi, \epsilon) = \sum_{\alpha=1}^N \left(\frac{1}{2} \epsilon^\alpha : \mathbf{C}^\alpha : \epsilon^\alpha \right) h^\alpha(\phi) = \frac{1}{2} \epsilon : \left(\sum_{\alpha=1}^N \mathbf{C}^\alpha h^\alpha(\phi) \right) : \epsilon = \frac{1}{2} \epsilon : \mathbf{C}_{VT}(\phi) : \epsilon.$$

Based on this form of f_{el} and the Equation (7.73), $\boldsymbol{\sigma}$ is easily seen to satisfy the important “compatibility” condition $\boldsymbol{\sigma}(\phi, \epsilon) = \frac{\partial f_{el}}{\partial \epsilon}(\phi, \epsilon)$ in Equation (7.70).

This can clearly be considered to be the mechanical analogue of the chemical model used e.g. in [52], whose underlying chemical free energy density can, as discussed above, be interpreted as a h^α -weighted interpolated of the phase-specific ones f^α under the assumption of equal phase-specific concentrations $\mathbf{c}^\alpha = \mathbf{c} \forall \alpha$. In addition, the total chemical potential $\boldsymbol{\mu}(\phi, \mathbf{c})$ is, similar to

⁴⁶Depending on the precise form of the implementation, this can nevertheless be somewhat tedious from a practical point of view.

the stress in Equation (7.73), given as the weighted average $\boldsymbol{\mu} = \sum_{\alpha} \mu^{\alpha}(\mathbf{c}) h^{\alpha}(\phi)$ of the phase-specific chemical potentials $\boldsymbol{\mu}^{\alpha}(\mathbf{c}) = \frac{\partial f^{\alpha}}{\partial \mathbf{c}^{\alpha}}(\mathbf{c})$.

In the second approach, the **Reuss-Sachs model** (also called **Steinbach-Apel model** based on their paper [69]), one assumes the equality of the phase-specific stresses with the total one, $\boldsymbol{\sigma}^{\alpha} = \boldsymbol{\sigma} \forall \alpha$, instead of that of the strains. This leads to the second equality in (7.72) being trivially satisfied, whereas, using $\boldsymbol{\epsilon}^{\alpha} = (\mathbf{C}^{\alpha})^{-1} : \boldsymbol{\sigma}^{\alpha} = (\mathbf{C}^{\alpha})^{-1} : \boldsymbol{\sigma} = \mathbf{S}^{\alpha} : \boldsymbol{\sigma}$, the first one imposes the restriction

$$\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^{\alpha} h^{\alpha}(\phi) = \left(\sum_{\alpha=1}^N \mathbf{S}^{\alpha} h^{\alpha}(\phi) \right) : \boldsymbol{\sigma} \quad (7.74)$$

and thus corresponds an effective material behavior where one interpolates the individual compliance tensors $\mathbf{S}^{\alpha} = (\mathbf{C}^{\alpha})^{-1}$ instead of the stiffness tensors. Inverting relation (7.74) leads to the explicit stress-strain relationship

$$\boldsymbol{\sigma}_{RS}(\phi, \boldsymbol{\epsilon}) = \left(\sum_{\alpha=1}^N \mathbf{S}^{\alpha} h^{\alpha}(\phi) \right)^{-1} : \boldsymbol{\epsilon} = \left(\sum_{\alpha=1}^N (\mathbf{C}^{\alpha})^{-1} h^{\alpha}(\phi) \right)^{-1} : \boldsymbol{\epsilon} = \mathbf{C}_{RS}(\phi) : \boldsymbol{\epsilon}, \quad (7.75)$$

with $\mathbf{C}_{RS}(\phi) = \left(\sum_{\alpha=1}^N (\mathbf{C}^{\alpha})^{-1} h^{\alpha}(\phi) \right)^{-1}$ and thus corresponds to a weighted harmonic interpolation of the stiffness tensors. In addition, using $\boldsymbol{\epsilon}^{\alpha} = \mathbf{S}^{\alpha} : \boldsymbol{\sigma}^{\alpha} = \mathbf{S}^{\alpha} : \boldsymbol{\sigma}$, the local energy contribution is given by

$$\begin{aligned} f_{el,RS}(\phi, \boldsymbol{\epsilon}) &= \sum_{\alpha=1}^N \left(\frac{1}{2} \boldsymbol{\epsilon}^{\alpha} : \mathbf{C}^{\alpha} : \boldsymbol{\epsilon}^{\alpha} \right) h^{\alpha}(\phi) = \frac{1}{2} \boldsymbol{\sigma} : \left(\sum_{\alpha=1}^N (\mathbf{S}^{\alpha} h^{\alpha}(\phi)) \right) : \boldsymbol{\sigma} \\ &= \frac{1}{2} \left(\left(\sum_{\alpha=1}^N \mathbf{S}^{\alpha} h^{\alpha}(\phi) \right)^{-1} : \boldsymbol{\epsilon} \right) : \left(\sum_{\alpha=1}^N (\mathbf{S}^{\alpha} h^{\alpha}(\phi)) \right) : \left(\left(\sum_{\alpha=1}^N \mathbf{S}^{\alpha} h^{\alpha}(\phi) \right)^{-1} : \boldsymbol{\epsilon} \right) \\ &= \frac{1}{2} \boldsymbol{\epsilon} : \left(\sum_{\alpha=1}^N (\mathbf{S}^{\alpha} h^{\alpha}(\phi)) \right)^{-1} : \boldsymbol{\epsilon} = \frac{1}{2} \boldsymbol{\epsilon} : \mathbf{C}_{RS}(\phi) : \boldsymbol{\epsilon} \end{aligned} \quad (7.76)$$

and is in particular again consistent with Equation (7.70).

In contrast to the Voigt-Taylor-model, this corresponds to a mechanical analogue (see also the discussion in [69]) of the KKS-type models considered in Section 7.1. The only slight difference is that in the chemical case one was only able to obtain equality of the reduced chemical potentials (resp. equality of the $\boldsymbol{\mu}^{\alpha} = \frac{\partial f^{\alpha}}{\partial \mathbf{c}^{\alpha}}$ only up to phase-specific constants). since there is no equivalence of the additional ‘‘internal’’ constraint $\sum_{i=1}^K c_i^{\alpha} = 1$ in the mechanical case, one can in fact enforce full equality of the phase-specific stresses.

Remark 117. In relation with the quantitative chemical models from Section 7.1, it is further interesting to consider how the above two models above compare from an energetic point of view. In fact, the Reuss-Sachs model corresponds to the minimizer of the local strain energy density $f_{el}(\phi, \boldsymbol{\epsilon})$ over all decompositions $\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^{\alpha} h^{\alpha}(\phi)$ into phaseinherent strains $(\boldsymbol{\epsilon}^{\alpha})_{1 \leq \alpha \leq N}$ as in Equation (7.72),

$$f_{el,RS}(\phi, \boldsymbol{\epsilon}) = \min_{\{\boldsymbol{\epsilon}^{\alpha} : \sum_{\alpha=1}^N \boldsymbol{\epsilon}^{\alpha} h^{\alpha}(\phi) = \boldsymbol{\epsilon}\}} \sum_{\alpha=1}^N \left(\frac{1}{2} \boldsymbol{\epsilon}^{\alpha} : \mathbf{C}^{\alpha} : \boldsymbol{\epsilon}^{\alpha} \right) h^{\alpha}(\phi). \quad (7.77)$$

This follows from the fact that for $0 \leq h^{\alpha}(\phi) \leq 1$, the problem (7.77) is a convex (and strictly convex for those $\boldsymbol{\epsilon}^{\alpha}$ with $h^{\alpha}(\phi) > 0$)⁴⁷ minimization problem and that, by a simple differentiation

⁴⁷Note that, even though convexity in each variable separately does not generally imply ‘‘global’’ convexity, it does so here due to the separated form of $f_{el,RS}$ in the $\boldsymbol{\epsilon}^{\alpha}$. In addition, the feasible set is clearly convex as any convex combination of phase-specific strains averaging to $\boldsymbol{\epsilon}$ will also do so.

of the Lagrangian,

$$L(\phi, (\epsilon^\alpha)_{1 \leq \alpha \leq N}, \sigma) = \sum_{\alpha=1}^N \left(\frac{1}{2} \epsilon^\alpha : \mathbf{C}^\alpha : \epsilon^\alpha \right) h^\alpha(\phi) - \sigma : \left(\sum_{\alpha=1}^N \epsilon^\alpha h^\alpha(\phi) - \epsilon \right),$$

the minimizers satisfy $(\mathbf{C}^\alpha : \epsilon^\alpha - \sigma) h^\alpha(\phi) = \mathbf{0}$, i.e. $\sigma = \mathbf{C}^\alpha : \epsilon^\alpha = \sigma^\alpha$ and thus the equality of the stresses has to hold for all α with $h^\alpha(\phi) \neq 0$ ⁴⁸.

Similarly, the stress-strain relationship for the Voigt-Taylor approach can be recovered from the minimization problem

$$f_{el,VT}^*(\phi, \sigma) = \min_{\{\sigma^\alpha : \sum_{\alpha=1}^N \sigma^\alpha h^\alpha(\phi) = \sigma\}} \sum_{\alpha=1}^N \left(\frac{1}{2} \sigma^\alpha : \mathbf{S}^\alpha : \sigma^\alpha \right) h^\alpha(\phi). \quad (7.78)$$

corresponding to a minimization of the weighted average of the complementary strain energy densities

$$(f_{el})^*(\phi, (\sigma^\alpha)_{1 \leq \alpha \leq N}) = \sum_{\alpha=1}^N (f_{el}^\alpha)^*(\sigma^\alpha) h^\alpha(\phi) = \sum_{\alpha=1}^N \left(\frac{1}{2} \sigma^\alpha : \mathbf{S}^\alpha : \sigma^\alpha \right) h^\alpha(\phi) \quad (7.79)$$

over all decompositions satisfying $\sigma = \sum_{\alpha=1}^N \sigma^\alpha h^\alpha(\phi)$. This is again a convex function of the $(\sigma^\alpha)_{1 \leq \alpha \leq N}$ and the critical points of the corresponding Lagrangian satisfy $(\mathbf{S}^\alpha : \sigma^\alpha - \epsilon) h^\alpha(\phi) = \mathbf{0}$ i.e. $\epsilon^\alpha = \mathbf{S}^\alpha : \sigma^\alpha = \epsilon$ for all α with $h^\alpha(\phi) > 0$.

Note that in both cases, even though one a priori only imposes one of the two interpolation properties in Equation (7.72) in this variational formulation, the other one is actually also trivially satisfied (the phase-specific stresses in the Reuss-Sachs model resp. the phase-specific strains in the Voigt-Taylor model all being equal by the FONC for the minimization problems).

An obvious consequence of the characterizations above is that (see e.g. also [51], [61]) one necessarily has the inequalities

$$f_{el,RS}(\phi, \epsilon) \leq f_{el,VT}(\phi, \epsilon) \quad \text{and} \quad f_{el,VT}^*(\phi, \sigma) \leq f_{el,RS}^*(\phi, \sigma)$$

as the Reuss-Sachs model is (according to Equation (7.77)) the one with the smallest possible elastic free energy density defined by (7.71) for a given strain and the interpolation property of the strains, whereas (according to Equation (7.78)), the Voigt-Taylor model has the smallest complementary strain energy density for a given stress and the interpolation property for the stresses.

◇

Remark 118. In addition, as pointed out in [61], by a well-known homogenization result, these two interpolations correspond to variational bounds on the effective behavior which would be expected from a real two-phase material.

It should be noted though that this result does **not** imply that the Voigt-Taylor model corresponds to an upper bound on the strain energy density $f_{el}(\phi, (\epsilon^\alpha)_{1 \leq \alpha \leq N})$ subject to the constraint of averaging to a given effective strain ϵ . In fact, it is easy to see that the effective free energy can be made arbitrarily large for a given ϵ while maintaining the interpolation properties above⁴⁹.

◇

⁴⁸The other ones being basically arbitrary, but without any contribution to $f_{el,RS}$. As previously for the chemical models, these phases still have a contribution to the derivative w.r.t. ϕ with the natural “definition” being the one setting $\sigma^\alpha = \sigma$, even if $h^\alpha(\phi) = 0$.

⁴⁹This is a **convex** maximization problem in the ϵ^α with an unbounded admissible set.

The Case With Eigenstrains

The approaches above can also be extended to the more general setting including eigenstrains $(\tilde{\epsilon}^\alpha)_{1 \leq \alpha \leq N}$ ⁵⁰, but now with, as above, $\sigma^\alpha = \mathbf{C}^\alpha : (\epsilon - \tilde{\epsilon}^\alpha)$ and $f_{el}(\phi, \epsilon)$ as in Equation (7.71). Assuming again, as in the Voigt-Taylor case, equality of the (total) strains then leads to the phase-specific stress $\sigma^\alpha = \mathbf{C}^\alpha : (\epsilon - \tilde{\epsilon}^\alpha)$ and therefore a total stress given by

$$\begin{aligned} \sigma_{VT}(\phi, \epsilon, \{\tilde{\epsilon}^\alpha\}_{1 \leq \alpha \leq N}) &= \sum_{\alpha=1}^N \mathbf{C}^\alpha : (\epsilon - \tilde{\epsilon}^\alpha) h^\alpha(\phi) \\ &= \left(\sum_{\alpha=1}^N \mathbf{C}^\alpha h^\alpha(\phi) \right) : \epsilon - \sum_{\alpha=1}^N \mathbf{C}^\alpha : \tilde{\epsilon}^\alpha h^\alpha(\phi) = \mathbf{C}_{VT}(\phi) : (\epsilon - \hat{\tilde{\epsilon}}), \end{aligned} \quad (7.80)$$

with $\hat{\tilde{\epsilon}}(\phi, \{\tilde{\epsilon}^\alpha\}_{1 \leq \alpha \leq N}) = \mathbf{C}_{VT}^{-1}(\phi) : \left(\sum_{\alpha=1}^N \mathbf{C}^\alpha : \tilde{\epsilon}^\alpha h^\alpha(\phi) \right)$ (see also e.g. [3]), i.e. an effective material behavior with the same stiffness but an eigenstrain given by a \mathbf{C}^α -weighted average of the intrinsic ones and in particular **not** satisfying the interpolation rule $\tilde{\epsilon}(\phi) = \sum_{\alpha=1}^N \tilde{\epsilon}^\alpha h^\alpha(\phi)$.

If one in contrast assumes, corresponding to the model of Khachaturyan [18], the equality of the elastic strains instead of the total strains and **defining** the total strain to be given by $\epsilon := \sum_{\alpha=1}^N (\epsilon_{el} + \tilde{\epsilon}^\alpha) h^\alpha(\phi)$, the interpolation

$$\sigma_K(\phi, \epsilon, \{\tilde{\epsilon}^\alpha\}_{1 \leq \alpha \leq N}) := \sum_{\alpha=1}^N \mathbf{C}^\alpha : \epsilon_{el} h^\alpha(\phi) = \mathbf{C}_{VT}(\phi) : \epsilon_{el} = \mathbf{C}_{VT}(\phi) : \left(\epsilon - \sum_{\alpha=1}^N \tilde{\epsilon}^\alpha h^\alpha(\phi) \right),$$

leads to an effective material behavior which satisfies $\sigma = \mathbf{C}_{VT}(\phi) : (\epsilon - \tilde{\epsilon})$ with the ‘‘more natural’’ definitions of ϵ and $\tilde{\epsilon} := \sum_{\alpha=1}^N \tilde{\epsilon}^\alpha h^\alpha(\phi)$.

Remark 119. Even though this model, in terms of the effective stress-strain relationship, may seem like the simpler - and thus in a way more elegant one (requiring only averaged quantities and no explicit phase-specific ones) - the higher simplicity is, from an energetic point of view, also its major drawback⁵¹ (for a similar discussion, see [3]). In fact, it is obvious that, if the consistency relation (7.70) is to hold, this stress-strain-relationship results from the elastic free energy density

$$f_{el,K}(\phi, \epsilon, \tilde{\epsilon}(\phi)) = \frac{1}{2} (\epsilon - \tilde{\epsilon}(\phi)) : \mathbf{C}_{VT}(\phi) : (\epsilon - \tilde{\epsilon}(\phi)), \quad (7.81)$$

and will generally not coincide with an average of the phase-specific free energy densities in (7.71) unless all eigenstrains $\tilde{\epsilon}^\alpha$ happen to be the same. \diamond

The situation for the Reuss-Sachs model is slightly simpler. Assuming the equality of the stresses, one has $\epsilon^\alpha = \tilde{\epsilon}^\alpha + \mathbf{S}^\alpha : \sigma$, $\alpha = 1, \dots, N$, and, again imposing $\epsilon = \sum_{\alpha=1}^N \epsilon^\alpha h^\alpha(\phi)$, a total strain given by

$$\epsilon = \sum_{\alpha=1}^N (\tilde{\epsilon}^\alpha + \mathbf{S}^\alpha : \sigma) h^\alpha(\phi) = \sum_{\alpha=1}^N \tilde{\epsilon}^\alpha h^\alpha(\phi) + \left(\sum_{\alpha=1}^N \mathbf{S}^\alpha h^\alpha(\phi) \right) : \sigma.$$

Defining $\tilde{\epsilon} = \sum_{\alpha=1}^N \tilde{\epsilon}^\alpha h^\alpha(\phi)$ and inverting this relationship shows that

$$\sigma_{RS} = \left(\sum_{\alpha=1}^N \mathbf{S}^\alpha h^\alpha(\phi) \right)^{-1} : (\epsilon - \tilde{\epsilon}) = \mathbf{C}_{RS}(\phi) : (\epsilon - \tilde{\epsilon}) \quad \text{with} \quad \mathbf{C}_{RS}(\phi) = \left(\sum_{\alpha=1}^N \mathbf{S}^\alpha h^\alpha(\phi) \right)^{-1}$$

and thus that the effective stiffness coincides with the expected expression and that the effective eigenstrain $\tilde{\epsilon}$ additionally obeys the standard interpolation property.

⁵⁰Possibly depending upon additional paramters such as the concentration or temperature.

⁵¹This is similar in spirit to Section 7.1, where the properties of the model can improve significantly if one does **not** assume the equality of the more convenient quantities (in this case the phase-specific concentrations c^α) and instead, at the cost of a more complex model, the equality of derived quantities which are much more natural to assume to be the same for all phases.

From this, it follows $\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha = \boldsymbol{S}^\alpha : \boldsymbol{C}_{RS}(\boldsymbol{\phi}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}})$. Reinserting this effective stress-strain relationship into the definition of the elastic energy, one obtains the total elastic energy density given by

$$f_{el,RS}(\boldsymbol{\phi}, \boldsymbol{\epsilon}, \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) = \frac{1}{2} \sum_{\alpha=1}^N \left(\boldsymbol{S}^\alpha : \boldsymbol{C}_{RS}(\boldsymbol{\phi}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}) \right) : \boldsymbol{C}^\alpha : \left(\boldsymbol{S}^\alpha : \boldsymbol{C}_{RS}(\boldsymbol{\phi}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}) \right) h^\alpha(\boldsymbol{\phi}).$$

Canceling \boldsymbol{C}^α with \boldsymbol{S}^α and observing that the only remaining phase-dependent term is given by $\sum_{\alpha=1}^N \boldsymbol{S}^\alpha h^\alpha(\boldsymbol{\phi}) = (\boldsymbol{C}_{RS}(\boldsymbol{\phi}))^{-1}$, this further simplifies to

$$f_{el,RS}(\boldsymbol{\phi}, \boldsymbol{\epsilon}, \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) = \frac{1}{2} (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) : \boldsymbol{C}_{RS}(\boldsymbol{\phi}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) \quad (7.82)$$

corresponding to simply replacing the total strain $\boldsymbol{\epsilon}$ in Equation (7.76) with the average elastic strain $\boldsymbol{\epsilon}_{el} = \boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}$.

The Driving Force

The various choices above imply different driving forces $\frac{\partial f_{el}}{\partial \boldsymbol{\phi}}$ in the phasefield Equation (6.73). In the simplest cases, i.e. in the Voigt-Taylor and Khachaturyan-model, the elastic free energy density is given as a fully explicit formula in terms of the phasefield values $\boldsymbol{\phi}$ and the total strains $\boldsymbol{\epsilon}$. For the Voigt-Taylor model, a straightforward differentiation with $\boldsymbol{\epsilon}^\beta = \boldsymbol{\epsilon}$ for all β shows that

$$\frac{\partial f_{el,VT}(\boldsymbol{\phi}, \boldsymbol{\epsilon}, \{\tilde{\boldsymbol{\epsilon}}^\beta\}_{1 \leq \beta \leq N})}{\partial \phi^\alpha} = \sum_{\beta=1}^N \frac{1}{2} (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\beta) : \boldsymbol{C}^\beta : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha} = \sum_{\beta=1}^N f_{el}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha},$$

i.e. the evolution of the phasefield is driven purely by the difference of the elastic free energy densities f^α . While reducing to the same model in the purely elastic case, the situation is different in the presence of eigenstrains for the Khachaturyan-model, since, based on equation (7.81) and $\frac{\partial f_{el,K}}{\partial \tilde{\boldsymbol{\epsilon}}} = -\boldsymbol{C}_{VT} : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) = -\boldsymbol{\sigma}_K$, one has

$$\frac{\partial f_{el,K}(\boldsymbol{\phi}, \boldsymbol{\epsilon}, \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}))}{\partial \phi^\alpha} = \sum_{\beta=1}^N \frac{1}{2} \left((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) : \boldsymbol{C}^\beta : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) - \boldsymbol{\sigma} : \tilde{\boldsymbol{\epsilon}}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha}$$

which in general can only be artificially related to a derivative involving the $f^\beta(\boldsymbol{\epsilon})$ underlying Equation (7.71) unless all eigenstrains coincide (in which case the distinction between the Voigt-Taylor and Khachaturyan model again becomes irrelevant).

The situation for the Reuss-Sachs-model is slightly more complicated as one has, using, analogous to the Khachaturyan-model, $\frac{\partial f_{el,RS}}{\partial \tilde{\boldsymbol{\epsilon}}} = -\boldsymbol{C}_{RS} : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) = -\boldsymbol{\sigma}_{RS}$,

$$\frac{\partial f_{el,RS}(\boldsymbol{\phi}, \boldsymbol{\epsilon}, \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}))}{\partial \phi^\alpha} = \frac{1}{2} (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) : \left(\frac{\partial \boldsymbol{C}_{RS}}{\partial \phi^\alpha}(\boldsymbol{\phi}) \right) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi})) - \boldsymbol{\sigma}_{RS} : \sum_{\beta=1}^N \tilde{\boldsymbol{\epsilon}}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}. \quad (7.83)$$

Here, in contrast to the two previous models, $\boldsymbol{C}_{RS}(\boldsymbol{\phi})$ as a function of $\boldsymbol{\phi}$ is only defined implicitly as the (pseudo-)inverse of $\boldsymbol{S}(\boldsymbol{\phi}) = \sum_{\beta=1}^N \boldsymbol{S}^\beta h^\beta(\boldsymbol{\phi})$. This is not a serious issue though as it is well-known that⁵²

$$\frac{\partial \boldsymbol{C}_{RS}(\boldsymbol{\phi})}{\partial \phi^\alpha} = -\boldsymbol{C}_{RS}(\boldsymbol{\phi}) : \left(\sum_{\beta=1}^N \boldsymbol{S}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} \right) : \boldsymbol{C}_{RS}(\boldsymbol{\phi}).$$

⁵²To see this, it suffices to differentiate the identity $\boldsymbol{C}_{RS}(\boldsymbol{\phi}) : \left(\sum_{\alpha=1}^N \boldsymbol{S}^\alpha h^\alpha(\boldsymbol{\phi}) \right) = \boldsymbol{I}^{(s)}$, which is legitimate under the standard coercivity assumption for the \boldsymbol{C}^α resp. \boldsymbol{S}^α (on the subspace of symmetric tensors).

While this expression in itself is somewhat cumbersome, it in fact allows for a major simplification in combination with Equation (7.83) as, by symmetry of \mathcal{C}_{RS} ,

$$\begin{aligned} \frac{1}{2}(\epsilon - \tilde{\epsilon}(\phi)) : \left(\frac{\partial \mathcal{C}_{RS}}{\partial \phi^\alpha}(\phi) \right) : (\epsilon - \tilde{\epsilon}(\phi)) &= -\frac{1}{2}(\epsilon - \tilde{\epsilon}(\phi)) : \mathcal{C}_{RS}(\phi) : \left(\sum_{\beta=1}^N \mathcal{S}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} \right) : \mathcal{C}_{RS}(\phi) : (\epsilon - \tilde{\epsilon}(\phi)) \\ &= -\frac{1}{2} \sigma_{RS} : \left(\sum_{\beta=1}^N \mathcal{S}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} \right) : \sigma_{RS}. \end{aligned}$$

Since by construction $\sigma_{RS} = \sigma^\beta$ for all β , this is the same as $-\frac{1}{2} \sum_{\beta=1}^N \sigma^\beta : \mathcal{S}^\beta : \sigma^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} = -\sum_{\beta=1}^N f^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$, thus leading to a total driving force given by

$$\frac{\partial f_{el,RS}}{\partial \phi^\alpha}(\phi, \epsilon, \{\tilde{\epsilon}^\beta\}_{1 \leq \beta \leq N}) = - \sum_{\beta=1}^N \left(f_{el}^\beta + \sigma(\phi, \epsilon) : \tilde{\epsilon}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha}. \quad (7.84)$$

Remark 120. As pointed out above (and already discussed in e.g. [69]), the Reuss-Sachs/Steinbach-Apel scheme is similar in spirit to the quantitative chemicals models from Section 7.1 with the ϵ^α taking the role of the phase-specific concentrations c^α and σ the role of the (reduced) chemical potential $\tilde{\mu}$. In contrast, the Voigt-Taylor model (and, up to a point, the Khatchuryan approach) is more similar in nature to the chemical model used in e.g. [52], which can be interpreted as assuming the equality of the phase-specific concentrations c^α with the average one. Even though these analogies are helpful for interpreting the mechanical models above, it needs to be realized that they are not “complete”. Compared to the seemingly quite satisfactory case previously considered in Section 7.1, it was observed early on (see e.g. [23]) that the results obtained using the Reuss-Sachs model are usually not as satisfactory as those from the chemical case and in particular leads to excess interface energy contributions except for very special mechanical settings.

There are several reasons why developing successful models in the mechanical case is more difficult than in the purely chemical one. A first one is that it is more common to have to deal with an additional “internal” forcing expressed through the eigenstrains in the mechanical than in the chemical case. The more crucial one lies in the non-local dependence of f_{el} on \mathbf{u} through its (symmetrized) gradient ϵ and the resulting nature of the equilibrium conditions themselves. The classical condition of the equilibration of the chemical potentials underlying the reasoning in Section 7.1 is in fact a highly restrictive one which is physically valid only for the very particular (but commonly used) settings of a **closed** system.

From a more mathematical point of view, this setting is related to the minimization of the free energy functional

$$\mathcal{F}_\epsilon(\phi, \mathbf{c}) = \int_{\Omega} \epsilon a(\phi, \nabla \phi) + \frac{1}{\epsilon} w(\phi) + f_{ch}(\phi, \mathbf{c}) \, d\mathbf{x} \quad (7.85)$$

under the implicit assumption that \mathcal{F}_ϵ is to be minimized with respect to \mathbf{c} under the additional constraint that the total amount of concentration of each component is to be maintained constant, $\int_{\Omega} c_i \, d\mathbf{x} = \text{const}$. Even though the form of the functional (7.85) does not lend itself naturally to the application of boundary conditions on \mathbf{c} ⁵³, the same effect can be enforced in terms of the gradient-flow approach to the minimization of \mathcal{F}_ϵ by “externally” postulating a conservative gradient-flow for \mathbf{c} as discussed in Section 6.1, where the conservation of the total mass is ensured by the locally conservative flow in combination with the “natural” isolating boundary conditions for the concentration.

⁵³As it is purely local in \mathbf{c} , there is a priori no reason for any minimizer in \mathbf{c} to be sufficiently regular (e.g. $\mathbf{c} \in \mathbf{H}^1(\Omega)$) to make sense of its boundary values such that this regularity would have to be enforced “artificially”.

In contrast, when modeling other physical situations (involving for example prescribed concentration values and/or chemical potential values on the outer boundary), the simple condition $\tilde{\boldsymbol{\mu}} = \text{const}$ will in general not hold anymore. Instead, the physically relevant equilibrium condition (in terms of the model in [52]) becomes

$$\begin{cases} -\nabla \cdot (\mathbf{L}_{ij}(\boldsymbol{\phi}, \mathbf{c}) \nabla \mu_j(\boldsymbol{\phi}, \mathbf{c})) = \mathbf{0} & \text{in } \Omega, \\ \text{boundary conditions on } \mathbf{c} \text{ resp. } \boldsymbol{\mu}. \end{cases} \quad (7.86)$$

On the one hand, it is clear (there for example not being any \mathbf{L}_{ij} in $\mathcal{F}_\epsilon(\boldsymbol{\phi}, \mathbf{c})$) that this is not easily connected to a minimization of the functional in Equation (7.85) unless one would treat this explicitly as a constrained optimization problem subject to Equation (7.86). On the other hand, as the condition $\tilde{\boldsymbol{\mu}} = \text{const}$ would have to be replaced with the (significantly) more complex equilibrium condition in Equation (7.86) corresponding to $-\text{div}(\boldsymbol{\sigma}(\boldsymbol{\phi}, \boldsymbol{\epsilon})) = \rho \mathbf{f}$ in the mechanical case. In terms of the diffuse interface approximation, it is easy to see that one would therefore also have to deal with the same issues as will be discussed in the mechanical case in the following sections.

The primary difference between the chemical and mechanical case therefore lies in the fact that, for the physically interesting examples, the latter one only very rarely allows for the simple solution $\boldsymbol{\sigma} = \text{const}$ ⁵⁴, whereas the analogous condition $\tilde{\boldsymbol{\mu}} = \text{const}$ is much more common in the chemical setting⁵⁵. \diamond

7.2.2 The Quantitative Model in the Two-Phase Case

A more recent modeling approach, first introduced in the two-phase setting by [51] and [64] (and later partially extended to the multiphase setting in [63], [62], [66] and also to situations involving plastic [35] and viscoelastic effects is based on the observation that, at an interface between two phases in the sharp interface setting, one usually neither has the equality of the strains nor that of the stresses. Instead, the spatial nature of the equilibrium conditions is also reflected in the well-known mechanical jump conditions on internal interfaces. These require on the one hand that the tangential components of the displacement gradient coincide on the interface, i.e. the jump $[[\nabla \mathbf{u}]]$ of $\nabla \mathbf{u}$ upon traversing the interface is oriented along the normal direction and the displacement gradients on both sides are therefore related by the condition $(\nabla \mathbf{u})^+ = (\nabla \mathbf{u})^- + \mathbf{a} \otimes \mathbf{n}$. For the strain tensor usually used within the small deformation setting, this obviously implies that $[[\boldsymbol{\epsilon}]] = (\mathbf{a} \otimes \mathbf{n})_S = \frac{1}{2}(\mathbf{a} \otimes \mathbf{n} + \mathbf{n} \otimes \mathbf{a})$. On the other hand, in contrast to $\nabla \mathbf{u}$, it is the normal components of the stress-tensor on both sides of the interface which need to coincide, i.e. $[[\boldsymbol{\sigma}]] \cdot \mathbf{n} = \mathbf{0}$.

Remark 121. The idea of treating the normal and tangential components was actually introduced slightly earlier by Durga et al. under an additional simplifying assumption in the two-dimensional setting [23] and later generalized by the same authors in [24] to the three-dimensional setting without this simplification. As will be discussed in Section 7.2.3, the models in [51] and [64] differ from - and in fact improve upon - the generalization in [24] in a small but important point. For this reason, despite the fundamental importance of the work [23] in the current setting for initiating this line of thought, the other model(s) considered below should be preferred over the

⁵⁴A notable exception being the purely elastic material chain, in which case the Reuss-Sachs approximation is in fact free of excess energy contributions, see e.g. [64].

⁵⁵There are also a significant number of works - for example in directional solidification - in which the models from Section 7.1 are applied in situations where this does not hold. These applications usually being set up (using e.g. moving window techniques) such that one never actually achieves equilibrium, the precise advantages and disadvantages of the models are significantly harder to judge as they are always based on the much more complex dynamic case. In particular, it is common to employ an artificial dynamic contribution in terms of an antitrapping current, further increasing the complexity of the situation.

one in [24]⁵⁶. ◇

Remark 122. It is useful to keep the reasoning behind these jump conditions in mind. The first one on $[[\nabla \mathbf{u}]]$ is essentially imposed a priori - at least within the standard variational approach - by demanding a global \mathbf{H}^1 -regularity of the solution. This implies that the trace $\gamma(\mathbf{u})$ of \mathbf{u} on a smooth hypersurface Γ can be uniquely defined as an element of $\mathbf{H}^{\frac{1}{2}}(\Gamma)$, i.e. denoting the limiting values of \mathbf{u} as one approaches this hypersurface from both sides by $\gamma^\pm(\mathbf{u})$, one has $\gamma^+(\mathbf{u}) = \gamma^-(\mathbf{u})$. The traces coinciding along the surface then implies the same for their tangential derivatives⁵⁷.

In contrast, the second condition on $\boldsymbol{\sigma}$ is only valid provided the right-hand side is regular enough and in particular (consistent with the intuitive reasoning) excludes the case of a “concentrated” surface force acting on Γ . Focusing on the simpler steady-state equilibrium conditions $-\operatorname{div}(\boldsymbol{\sigma}) = \mathbf{f}$, \mathbf{f} being in e.g. $\mathbf{L}^2(\Omega)$ implies (by the very definition) that $\boldsymbol{\sigma} \in \mathbf{H}_{\operatorname{div}}(\Omega)$ ⁵⁸. Similar to the trace of functions in $\mathbf{H}^1(\Omega)$, it is a standard fact (see e.g. [34], [43] or [48]) that one can define a unique **normal** trace in $\mathbf{H}^{-\frac{1}{2}}(\Gamma)$ on any (sufficiently smooth) hypersurface for such functions⁵⁹, and that therefore the normal components of $\boldsymbol{\sigma}$ have to coincide as one approaches Γ from both sides. ◇

More elaborate mechanical models taking these sharp interface jump conditions at a two-phase interface into account have first been proposed by Mosler et al. [51] within an energy-based setting (a priori relying upon the jump conditions for the strains alone) and by Durga et al. [24] and Schneider et al. [64] within a more direct approach based on the jump conditions for both the strains and the stresses. In contrast to the work by [51], these models do not rely on a jump vector \mathbf{a} (as also used below), but instead use a basis transformation into an appropriate coordinate system in which they then enforce the equality of the appropriate stress- and strain-entries. Yet another approach is taken in [74], which avoids the use of the basis transformation by using a description in terms of projectors \mathcal{N} and \mathcal{T} onto the “normal” and “tangent” subspaces of the symmetric second-order strain and stress tensors.

Given that all models are essentially based on the same sharp-interface jump conditions, it is not surprising that, from a modeling point of view, they are in fact all different representations of the same physical model. Nevertheless, the choice of representation clearly can have a significant impact both in terms of the computational complexity of the models as well as the effort required for their mathematical description and practical implementation. For this reason, after a quick outline of the mechanical model and its additional contributions to the phasefield equation - here using a description in terms of a jump vector \mathbf{a} - the links between the various models will be discussed in Section 7.2.3 in a little more detail.

Further discussions, extensions and applications (often within the multiphase setting to be discussed in the next section) can e.g. be found in [61], [63], [74], [62], [35], [66], [4] and [5].

The Mechanical Model

The underlying idea is to increase the accuracy of the phasefield model (with a small but non-zero interface width) by directly incorporating the jump conditions from the sharp-interface setting. Following [51] (see also [61] for a more detailed discussion), one can therefore impose

⁵⁶Nevertheless, as shown in [23], their model does, among other things, in fact improve significantly upon the Voigt-Taylor or Reuss-Sachs model in terms of interfacial excess energies.

⁵⁷Then in an even weaker sense as elements of $\mathbf{H}^{-\frac{1}{2}}(\Gamma)$.

⁵⁸ $\mathbf{H}_{\operatorname{div}}(\Omega)$ being defined by $\mathbf{H}_{\operatorname{div}}(\Omega) := \{\boldsymbol{\sigma} \in \mathbf{L}^2(\Omega) : \operatorname{div}(\boldsymbol{\sigma}) \in \mathbf{L}^2(\Omega)\}$.

⁵⁹The essential idea is based upon “extending” Gauss’s divergence theorem $\int_\omega \operatorname{div}(\mathbf{g}) \, d\mathbf{x} = \int_{\partial\omega} \mathbf{g} \cdot \mathbf{n} \, ds$ for smooth functions by appropriately **defining** the (a priori meaningless in the non-smooth case) values of \mathbf{g} such that this integral equality holds.

If the divergence of $\boldsymbol{\sigma}$ is for example only in $\mathbf{H}^{-1}(\Omega)$ (e.g. if surface forces are present), the definition of an appropriate trace is a much trickier issue (see e.g. [48]).

the condition

$$\llbracket \mathbf{F} \rrbracket = \mathbf{F}^2 - \mathbf{F}^1 = \mathbf{a} \otimes \mathbf{n} \quad \text{resp.} \quad \llbracket \boldsymbol{\epsilon} \rrbracket = \boldsymbol{\epsilon}^2 - \boldsymbol{\epsilon}^1 = (\mathbf{a} \otimes \mathbf{n})_S \quad (7.87)$$

on the displacement gradients $\mathbf{F} = \nabla \mathbf{u}$ resp. its symmetrization in the small deformation setting. Combined with the interpolation requirement $\boldsymbol{\epsilon} = \boldsymbol{\epsilon}^1 h^2(\phi) + \boldsymbol{\epsilon}^2 h^1(\phi)$, where, as in Section 6.2, $h^1(\phi) = h(\phi)$, $h^2 = h(1 - \phi) = 1 - h(\phi)$, this results in an effective strain given by $\boldsymbol{\epsilon} = \boldsymbol{\epsilon}^1 h^1(\phi) + (\boldsymbol{\epsilon}^1 + (\mathbf{a} \otimes \mathbf{n})_S) h^2(\phi) = \boldsymbol{\epsilon}^1 + h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S$, which can be solved for the strains $(\boldsymbol{\epsilon}^1, \boldsymbol{\epsilon}^2)$ as a function of $\boldsymbol{\epsilon}$ and the (as of yet arbitrary) jump vector \mathbf{a} to give

$$\boldsymbol{\epsilon}^1(\boldsymbol{\epsilon}, \mathbf{a}) = \boldsymbol{\epsilon} - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S \quad \text{and} \quad \boldsymbol{\epsilon}^2(\boldsymbol{\epsilon}, \mathbf{a}) = \boldsymbol{\epsilon}^1 + (\mathbf{a} \otimes \mathbf{n})_S = \boldsymbol{\epsilon} + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S. \quad (7.88)$$

Inserting these expressions into the material law $\boldsymbol{\sigma}^\alpha = \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ and the jump condition for $\boldsymbol{\sigma}$, one obtains the additional condition

$$\begin{aligned} \llbracket \boldsymbol{\sigma} \rrbracket \cdot \mathbf{n} &= (\mathbf{C}^2 : (\boldsymbol{\epsilon}^2 - \tilde{\boldsymbol{\epsilon}}^2) - \mathbf{C}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1)) \cdot \mathbf{n} \\ &= \left(\mathbf{C}^2 : ((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S) - \mathbf{C}^1 : ((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S) \right) \cdot \mathbf{n} \stackrel{!}{=} \mathbf{0}. \end{aligned}$$

Given the subsymmetry $\mathcal{C}_{ijkl}^\alpha = \mathcal{C}_{ijlk}^\alpha$ of the stiffness-tensor⁶⁰, this can in fact be simplified to the equation

$$\left(\mathbf{n} \cdot \left(h^2(\phi) \mathbf{C}^1 + h^1(\phi) \mathbf{C}^2 \right) \cdot \mathbf{n} \right) \cdot \mathbf{a} \stackrel{!}{=} \left(\mathbf{C}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - \mathbf{C}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) \right) \cdot \mathbf{n}.$$

This is a 3×3 linear system allowing for the determination of \mathbf{a} in terms of $\boldsymbol{\epsilon}$ as

$$\begin{aligned} \mathbf{a} &= \left(\mathbf{n} \cdot \left(h^2(\phi) \mathbf{C}^1 + h^1(\phi) \mathbf{C}^2(\phi) \right) \cdot \mathbf{n} \right)^{-1} \cdot \left(\mathbf{C}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - \mathbf{C}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) \right) \cdot \mathbf{n} \\ &= \left(\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n} \right)^{-1} \cdot \left(\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2 \right) \cdot \mathbf{n}, \end{aligned} \quad (7.89)$$

where, adopting the notation in [74], $\bar{\mathcal{C}}^{12}(\phi)$ denotes the ‘‘anti-arithmetic’’ average

$$\bar{\mathcal{C}}^{12}(\phi) := h^2(\phi) \mathbf{C}^1 + h^1(\phi) \mathbf{C}^2(\phi) = (1 - h(\phi)) \mathbf{C}^1 + h(\phi) \mathbf{C}^2 \quad (7.90)$$

and the $\boldsymbol{\Sigma}^\alpha := \mathbf{C}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)$ correspond to the phase-specific stresses predicted by the Voigt-Taylor model.

Remark 123. The solvability of the system for \mathbf{a} (resp. the invertibility of $\bar{\mathcal{C}}^{12}$ under the standard assumptions on the \mathbf{C}^α can e.g. be seen from the fact that the bilinear form $\frac{1}{2} \boldsymbol{\epsilon} : \mathbf{C}^\alpha : \boldsymbol{\epsilon}$ is strictly convex in $\boldsymbol{\epsilon}$, which implies rank-one convexity. Together with the subsymmetries of the \mathbf{C}^α and as $\bar{\mathcal{C}}^{12}$ is a convex combination (for $0 \leq h(\phi) \leq 1$) of \mathbf{C}^1 and \mathbf{C}^2 , it follows that $\mathbf{a} \mapsto \mathbf{a} \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})$ is strictly convex and $\bar{\mathcal{C}}^{12}(\phi)$ a positive definite matrix. \diamond

Remark 124. Depending on the way the stiffness-tensors are arranged in memory, it may be advantageous to use any of the (equivalent given the subsymmetries) expressions

$$n_i \mathcal{C}_{ijkl} n_l = n_j \mathcal{C}_{ijkl} n_l = n_j \mathcal{C}_{ijkl} n_k = n_i \mathcal{C}_{ijkl} n_k.$$

Note that none of these is the same as $\mathbf{C} : (\mathbf{n} \otimes \mathbf{n}) = (\mathbf{n} \otimes \mathbf{n}) : \mathbf{C}$ though.

With respect to the commonly used Voigt-notation, it should also be noted that, as this formula

⁶⁰By the right subsymmetric of \mathbf{C} , $\mathcal{C}_{ijkl} = \mathcal{C}_{ijlk}$, applying \mathbf{C} to a non-symmetric tensor \mathbf{t} leads to the same result as applying it to the symmetrized tensor $\mathbf{t}_S = \frac{1}{2}(\mathbf{t} + \mathbf{t}^T)$ since (using the Einstein summation convention)

$$\mathbf{C} : \mathbf{t}_S = \frac{1}{2}(\mathcal{C}_{ijkl} t_{kl} + \mathcal{C}_{ijlk} t_{lk}) = \frac{1}{2}(\mathcal{C}_{ijkl} t_{kl} + \mathcal{C}_{ijlk} t_{kl}) \stackrel{\mathcal{C}_{ijkl} = \mathcal{C}_{ijlk}}{=} \mathcal{C}_{ijkl} \frac{1}{2}(t_{kl} + t_{kl}) = \mathbf{C} : \mathbf{t}.$$

involves two separate contractions (and not a double-contraction with a symmetric tensor in a “reduced” vector format) there are no additional prefactors involved in the evaluation of $\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n}$. \diamond

With \mathbf{a} determined by Equation (7.89) and with $\boldsymbol{\sigma}^\alpha = \mathcal{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ together with Equation (7.88), one can calculate the phase-specific stresses (again using the right subsymmetry of the stiffness) to be given by

$$\begin{aligned}\boldsymbol{\sigma}^1 &= \mathcal{C}^1 \left(\boldsymbol{\epsilon} - h^2(\phi) (\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^1 \right) = \boldsymbol{\Sigma}^1 - h^2(\phi) (\mathcal{C}^1 \cdot \mathbf{n}) \cdot \mathbf{a} \\ &= \boldsymbol{\Sigma}^1 - h^2(\phi) (\mathcal{C}^1 \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) \cdot \mathbf{n} \right)\end{aligned}$$

and similarly

$$\boldsymbol{\sigma}^2 = \boldsymbol{\Sigma}^2 + h^1(\phi) (\mathcal{C}^2 \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) \cdot \mathbf{n} \right).$$

Finally, together with the definition of $\boldsymbol{\sigma}$ as the weighted average of $\boldsymbol{\sigma}^1$ and $\boldsymbol{\sigma}^2$, this leads to the effective material behavior given by

$$\begin{aligned}\boldsymbol{\sigma} &= h^1(\phi) \boldsymbol{\sigma}^1 + h^2(\phi) \boldsymbol{\sigma}^2 = h^1(\phi) \boldsymbol{\Sigma}^1 + h^2(\phi) \boldsymbol{\Sigma}^2 - h^1(\phi) h^2(\phi) \left((\mathcal{C}^1 - \mathcal{C}^2) \cdot \mathbf{n} \right) \cdot \mathbf{a} \\ &= \boldsymbol{\sigma}_{VT}(\boldsymbol{\epsilon}, \phi) - \left((\mathcal{C}^1 - \mathcal{C}^2) \cdot \mathbf{n} \right) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) \cdot \mathbf{n} \right),\end{aligned}\tag{7.91}$$

in terms of the various auxiliary quantities introduced above (and with $\boldsymbol{\sigma}_{VT}(\boldsymbol{\epsilon}, \phi)$ corresponding to the stress from the Voigt-Taylor model in equations (7.73) resp. (7.80)), or, after reinserting the definitions, the fully explicit formula

$$\begin{aligned}\boldsymbol{\sigma} &= h^1(\phi) \mathcal{C}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) + h^2(\phi) \mathcal{C}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) \\ &\quad - h^1(\phi) h^2(\phi) \left((\mathcal{C}^1 - \mathcal{C}^2) \cdot \mathbf{n} \right) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot \left(\mathbf{n} \cdot \left(\mathcal{C}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - \mathcal{C}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) \right) \right)\end{aligned}\tag{7.92}$$

for $\boldsymbol{\sigma}$ in terms of $\boldsymbol{\epsilon}$, \mathbf{n} and ϕ .

Remark 125. The formulas above simplify somewhat when there are no eigenstrains, i.e. when $\tilde{\boldsymbol{\epsilon}}^1 = \tilde{\boldsymbol{\epsilon}}^2 = \mathbf{0}$. It is easy to see that one then obtains

$$\boldsymbol{\sigma}^1 = \mathcal{C}^1 : \boldsymbol{\epsilon} - h^2(\phi) (\mathbf{n} \cdot \mathcal{C}^1) : \left((\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot (\mathbf{n} \cdot (\mathcal{C}^1 - \mathcal{C}^2) : \boldsymbol{\epsilon}) \right)$$

and

$$\boldsymbol{\sigma}^2 = \mathcal{C}^2 : \boldsymbol{\epsilon} + h^1(\phi) (\mathbf{n} \cdot \mathcal{C}^2) : \left((\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot (\mathbf{n} \cdot (\mathcal{C}^1 - \mathcal{C}^2) : \boldsymbol{\epsilon}) \right)$$

and thus, using the abbreviation $\mathcal{C}_{VT}(\phi) := h^1(\phi) \mathcal{C}^1 + h^2(\phi) \mathcal{C}^2 = h(\phi) \mathcal{C}^1 + (1 - h(\phi)) \mathcal{C}^2$ for the “arithmetic” Voigt-Taylor-type average of the stiffnesses and the various symmetries,

$$\boldsymbol{\sigma} = \left(\mathcal{C}_{VT}(\phi) - h^1(\phi) h^2(\phi) \left((\mathcal{C}^1 - \mathcal{C}^2) \cdot \mathbf{n} \right) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot (\mathbf{n} \cdot (\mathcal{C}^1 - \mathcal{C}^2)) \right) : \boldsymbol{\epsilon},\tag{7.93}$$

i.e. an effective behavior corresponding to the stiffness-tensor

$$\mathcal{C}_{eff}(\phi, \mathbf{n}) = \mathcal{C}_{VT}(\phi) - h^1(\phi) h^2(\phi) \left((\mathcal{C}^1 - \mathcal{C}^2) \cdot \mathbf{n} \right) \cdot (\mathbf{n} \cdot \bar{\mathcal{C}}^{12}(\phi) \cdot \mathbf{n})^{-1} \cdot (\mathbf{n} \cdot (\mathcal{C}^1 - \mathcal{C}^2)).\tag{7.94}$$

This effective stiffness can also be very useful for numerical purposes when using an implicit time-integration scheme or when making a quasi-steady state assumption for the elastic fields. As is well known, applying any of the standard Krylov solvers for the determination of the displacement field \mathbf{u} requires, except for the evaluation of the initial residual, only the ability to evaluate the system response to increments in \mathbf{u} , and thus here the ability to evaluate $\delta\boldsymbol{\sigma} = \boldsymbol{\sigma}(\boldsymbol{\epsilon}(\mathbf{u} + \delta\mathbf{u})) - \boldsymbol{\sigma}(\boldsymbol{\epsilon}(\mathbf{u}))$. As long as any eigenstrains considered do not have any additional dependence⁶¹ on $\boldsymbol{\epsilon}$, the stress-contributions due to the $\tilde{\boldsymbol{\epsilon}}^\alpha$ will cancel in $\delta\boldsymbol{\sigma}$ and the ‘‘homogeneous’’ version (7.93) and thus the effective stiffness $\mathcal{C}_{eff}(\phi, \mathbf{n})$ in Equation (7.94) is also the relevant one for evaluating the effect of increments in \mathbf{u} , $\delta\boldsymbol{\sigma} = \mathcal{C}_{eff}(\phi, \mathbf{n}) : \boldsymbol{\epsilon}(\delta\mathbf{u})$.

As $\mathcal{C}_{eff}(\phi)$ can in principle be precalculated and stored once at the beginning of each time-step, it is clear that the evaluation of this stress-increment using this prestored stiffness can be significantly cheaper than the procedure above (i.e. first determining the Voigt-Taylor stress predictions, then solving for \mathbf{a} and finally correcting $\boldsymbol{\sigma}_{VT}$). \diamond

Remark 126. For isotropic materials, the formula above can actually be evaluated analytically. Based on $\bar{\mathcal{C}}^{12} = \bar{\lambda}^{12}\mathbf{I} \otimes \mathbf{I} + 2\bar{\mu}^{12}\mathbf{I}^{(s)}$ where $\mathbf{I}^{(s)} = \frac{1}{2}(\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk})\mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_l$ is the symmetric fourth order unit tensor and $\bar{\lambda}^{12}$ and $\bar{\mu}^{12}$ the ‘‘anti-arithmetic’’ averages of λ and μ , one has $\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n} = \bar{\lambda}^{12}\mathbf{n} \otimes \mathbf{n} + 2\bar{\mu}^{12}n_i(\mathbf{I}^{(s)})_{ijkl}n_l\mathbf{e}_j \otimes \mathbf{e}_k$, or, as

$$n_i(\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk})n_l\mathbf{e}_j \otimes \mathbf{e}_k = (n_in_j\mathbf{e}_j \otimes \mathbf{e}_i + n_in_i\mathbf{e}_j \otimes \mathbf{e}_j) = (\mathbf{n} \otimes \mathbf{n} + \|\mathbf{n}\|^2\mathbf{I}) = (\mathbf{n} \otimes \mathbf{n} + \mathbf{I}),$$

the simple expression

$$\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n} = (\bar{\lambda}^{12} + \bar{\mu}^{12})\mathbf{n} \otimes \mathbf{n} + \bar{\mu}^{12}\mathbf{I}, \quad \text{where} \quad \bar{\lambda}^{12} = h^2\lambda^1 + h^1\lambda^2, \quad \bar{\mu}^{12} = h^2\mu^1 + h^1\mu^2$$

for $\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n}$ in terms of the Lamé-parameters λ and μ . Since this is a rank-one perturbation of the (scaled) identity, the inverse of this tensor is easily determined using the Sherman-Morrisson-Woodsbury-formula to be given by⁶²

$$(\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n})^{-1} = (\bar{\mu}^{12})^{-1} \left(\mathbf{I} - \frac{\bar{\lambda}^{12} + \bar{\mu}^{12}}{\bar{\mu}^{12}} \mathbf{n} \otimes \mathbf{n} \right) = (\bar{\mu}^{12})^{-1} \left(\mathbf{I} - \frac{\bar{\lambda}^{12} + \bar{\mu}^{12}}{\bar{\lambda}^{12} + 2\bar{\mu}^{12}} \mathbf{n} \otimes \mathbf{n} \right). \quad (7.95)$$

\diamond

The stress-strain relationship (7.92) above can alternatively - instead of defining it based on the weighted average of the two $\boldsymbol{\sigma}^\alpha = \mathcal{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$, $\alpha = 1, 2$ - also be obtained as the derivative $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$ of the elastic free energy density

$$\begin{aligned} f_{el} = & \frac{1}{2} \left(\boldsymbol{\epsilon} - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^1 \right) : \mathcal{C}^1 : \left(\boldsymbol{\epsilon} - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^1 \right) h^1(\phi) \\ & + \frac{1}{2} \left(\boldsymbol{\epsilon} + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^2 \right) : \mathcal{C}^2 : \left(\boldsymbol{\epsilon} + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^2 \right) h^2(\phi). \end{aligned} \quad (7.96)$$

This correspond to Equation (7.71) with $\boldsymbol{\epsilon}^\alpha$, $\alpha = 1, 2$ chosen as in Equation (7.88) provided the vector \mathbf{a} is, as above, fixed by the continuity condition $\boldsymbol{\sigma}^1 \cdot \mathbf{n} = \boldsymbol{\sigma}^2 \cdot \mathbf{n}$ on the normal stress components. In fact, even though \mathbf{a} depends on $\boldsymbol{\epsilon}$ through the defining normal continuity condition,

⁶¹Two such examples will quickly be considered below.

⁶²Alternatively, as in [74], one can also rewrite the formula for $\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n}$ as

$$\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n} = (\bar{\lambda}^{12} + 2\bar{\mu}^{12})\mathbf{n} \otimes \mathbf{n} + \bar{\mu}^{12}(\mathbf{I} - \mathbf{n} \otimes \mathbf{n})$$

where $\mathbf{n} \otimes \mathbf{n}$ and $\mathbf{I} - \mathbf{n} \otimes \mathbf{n}$ are the two orthogonal projection operators onto the normal subspace $\text{Span}(\{\mathbf{n}\})$ and the tangential one, $\text{Ker}(\mathbf{n})$. By orthogonality both terms can be inverted separately leading to

$$(\mathbf{n} \cdot \bar{\mathcal{C}}^{12} \cdot \mathbf{n})^{-1} = \frac{1}{\bar{\lambda}^{12} + 2\bar{\mu}^{12}} \mathbf{n} \otimes \mathbf{n} + \frac{1}{\bar{\mu}^{12}} (\mathbf{I} - \mathbf{n} \otimes \mathbf{n})$$

which is easily seen to be equivalent to the expression in Equation (7.95).

differentiating the free energy density w.r.t. ϵ and making use of the right subsymmetries of the \mathcal{C}^α leads to

$$\begin{aligned}\frac{\partial f_{el}}{\partial \epsilon} &= h^1 \mathcal{C}^1(\epsilon - h^2(\mathbf{a} \otimes \mathbf{n})_S) : (\mathbf{I}^{(s)} - h^2 \mathbf{n} \otimes \frac{\partial \mathbf{a}}{\partial \epsilon}) + h^2 \mathcal{C}^2(\epsilon + h^1(\mathbf{a} \otimes \mathbf{n})_S) : (\mathbf{I}^{(s)} + h^1 \mathbf{n} \otimes \frac{\partial \mathbf{a}}{\partial \epsilon}) \\ &= h^1 \boldsymbol{\sigma}^1 : (\mathbf{I}^{(s)} - h^2 \mathbf{n} \otimes \frac{\partial \mathbf{a}}{\partial \epsilon}) + h^2 : \boldsymbol{\sigma}^2 : (\mathbf{I}^{(s)} + h^1 \mathbf{n} \otimes \frac{\partial \mathbf{a}}{\partial \epsilon}).\end{aligned}$$

While the two terms $h^1 \boldsymbol{\sigma}^1 + h^2 \boldsymbol{\sigma}^2$ give rise to the desired effective stress $\boldsymbol{\sigma}$, the remaining two terms can be combined to $h^1 h^2 ((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}}{\partial \epsilon}$ and thus simply drop out by the imposed continuity of $\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}$. The two-phase model above therefore indeed does satisfy the important relation $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \epsilon}$ in Equation (7.70).

Remark 127. The crucial point here is that there is no additional contribution due to $\frac{\partial \mathbf{a}}{\partial \epsilon}$. The same will also be seen to hold with respect to the dependence of \mathbf{a} on the other parameters ϕ and \mathbf{n} and is in fact to be expected from the variational characterization by [51] to be recalled below. As will be seen in Section 7.2.4, the same property does unfortunately not hold for all of the multiphase generalizations which have been proposed, and then leads to significant complications as compared to the very favorable two-phase setting. \diamond

The Mechanical Driving Force

In contrast to the simpler models using an interpolation in terms of the local values of ϕ alone, the mechanical free energy density f_{el} now depends not only on the values of ϕ and ϵ alone, but has an additional dependence on the gradients of ϕ due to the heavy use of the normal vector \mathbf{n} for constructing the phase-specific stress- and strain-fields. Due to this additional dependence, there are now two different contributions to the mechanical driving for the phasefield model, namely the one due to the derivative w.r.t. ϕ itself and an additional divergence-type contribution $-\nabla \cdot (\frac{\partial f_{el}}{\partial \nabla \phi})$ due to the dependence on the orientation of the interface through $\mathbf{n} = -\nabla \phi$ ⁶³.

Remark 128. The appearance of an additional contribution to the phasefield equation due to the dependence on \mathbf{n} was already clearly stated in [51] but is in contrast (erroneously) claimed in [64] (and later corrected e.g. in [74]) not to be the case⁶⁴. \diamond

By the formula (7.89) for \mathbf{a} , this vector obviously depends both on ϕ and, through \mathbf{n} , on $\nabla \phi$, $\mathbf{a} = \mathbf{a}(\phi, \nabla \phi, \epsilon)$. Making use of $\frac{\partial f_{el}}{\partial \epsilon^\alpha} = \mathcal{C}^\alpha : (\epsilon^\alpha - \tilde{\epsilon}^\alpha) = \boldsymbol{\sigma}^\alpha$, differentiating the free energy density with respect to ϕ leads to

$$\begin{aligned}\frac{\partial f_{el}}{\partial \phi} &= f_{el}^1 \frac{\partial h^1}{\partial \phi} + f_{el}^2 \frac{\partial h^2}{\partial \phi} + \boldsymbol{\sigma}^1 : \frac{\partial \epsilon^1}{\partial \phi} h^1 + \boldsymbol{\sigma}^2 : \frac{\partial \epsilon^2}{\partial \phi} h^2 \\ &= f_{el}^1 \frac{\partial h^1}{\partial \phi} + f_{el}^2 \frac{\partial h^2}{\partial \phi} + \boldsymbol{\sigma}^1 : \left(-\frac{\partial h^2}{\partial \phi} (\mathbf{a} \otimes \mathbf{n})_S - h^2 \left(\frac{\partial \mathbf{a}}{\partial \phi} \otimes \mathbf{n} \right)_S \right) h^1 + \boldsymbol{\sigma}^2 : \left(\frac{\partial h^1}{\partial \phi} (\mathbf{a} \otimes \mathbf{n})_S + h^1 \left(\frac{\partial \mathbf{a}}{\partial \phi} \otimes \mathbf{n} \right)_S \right) h^2.\end{aligned}$$

As indicated in Remark 127, the dependence of \mathbf{a} on ϕ actually again drops out since the two terms containing $\frac{\partial \mathbf{a}}{\partial \phi}$ can be combined to give $h^1 h^2 (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \left(\frac{\partial \mathbf{a}}{\partial \phi} \otimes \mathbf{n} \right)_S = h^1 h^2 \left((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n} \right) \cdot \frac{\partial \mathbf{a}}{\partial \phi} = \mathbf{0}$ by the symmetry of $\boldsymbol{\sigma}$, and the normal stress-jump $\llbracket \boldsymbol{\sigma} \rrbracket \cdot \mathbf{n}$ vanishes by construction. The remaining terms can then, using $\frac{\partial h^2}{\partial \phi} = -\frac{\partial h^1}{\partial \phi}$ and $-h^1 \boldsymbol{\sigma}^1 \frac{\partial h^2}{\partial \phi} + h^2 \boldsymbol{\sigma}^2 \frac{\partial h^1}{\partial \phi} = (h^1 \boldsymbol{\sigma}^1 + h^2 \boldsymbol{\sigma}^2) \frac{\partial h^1}{\partial \phi} = \boldsymbol{\sigma} \frac{\partial h^1}{\partial \phi}$, be summarized to

$$\frac{\partial f_{el}}{\partial \phi} = (f_{el}^1 - f_{el}^2) \frac{\partial h^1}{\partial \phi} + \boldsymbol{\sigma} : (\mathbf{a} \otimes \mathbf{n})_S \frac{\partial h^1}{\partial \phi} = (f_{el}^1 - f_{el}^2) \frac{\partial h^1}{\partial \phi} + (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{a} \frac{\partial h^1}{\partial \phi}, \quad (7.97)$$

⁶³Note that in the two-phase setting, the assumption that \mathbf{n} is a function of $\nabla \phi$ alone is very natural. This can potentially be different in the multi-phase setting as when e.g. defining the normal(s) in terms of the $\mathbf{q}^{\alpha\beta}$ from Section 6.1 (making the normal(s) a function of both $\nabla \phi$ and ϕ) and thus leading to an additional contribution to $\frac{\partial f_{el}}{\partial \phi^\alpha}$ itself as well.

⁶⁴The error in [64] is hidden in the claim that “the partial derivatives are evaluated at constant $\boldsymbol{\epsilon}_B^\alpha = (\boldsymbol{\epsilon}_n, \boldsymbol{\epsilon}_t)$ ”, where $\boldsymbol{\epsilon}_n$ and $\boldsymbol{\epsilon}_t$ are the normal and tangential strain components.

or, with $(\mathbf{a} \otimes \mathbf{n})_S = \boldsymbol{\epsilon}^2 - \boldsymbol{\epsilon}^1$, to the alternative expression

$$\frac{\partial f_{el}}{\partial \phi} = \left((f_{el}^1 - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^1) - (f_{el}^2 - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^2) \right) \frac{\partial h^1}{\partial \phi} \quad (7.98)$$

using only the phase-specific strains.

Proceeding similarly for the derivative with respect to $\nabla \phi$, one has⁶⁵

$$\begin{aligned} \frac{\partial f_{el}}{\partial \nabla \phi} &= \frac{\partial}{\partial \nabla \phi} \left(-h^1 h^2 (\boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S + h^1 h^2 \boldsymbol{\sigma}^2 : (\mathbf{a} \otimes \mathbf{n})_S) \right) \\ &= h^1 h^2 \frac{\partial}{\partial \mathbf{n}} \left(((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}) \cdot \mathbf{n} \right) \cdot \frac{\partial \mathbf{n}}{\partial \nabla \phi} + h^1 h^2 \frac{\partial}{\partial \mathbf{a}} \left(((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n}) \cdot \mathbf{a} \right) \cdot \frac{\partial \mathbf{a}}{\partial \nabla \phi} \\ &= h^1 h^2 ((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}) \cdot \frac{\partial \mathbf{n}}{\partial \nabla \phi}, \end{aligned} \quad (7.99)$$

where the term $\frac{\partial \mathbf{a}}{\partial \nabla \phi}$ drops out by the continuity of the normal stresses. Using

$$\frac{\partial \mathbf{n}}{\partial \nabla \phi} = \frac{\partial}{\partial \nabla \phi} \left(-\frac{\nabla \phi}{|\nabla \phi|} \right) = -\frac{|\nabla \phi| \mathbf{I} - \nabla \phi \otimes \frac{\nabla \phi}{|\nabla \phi|}}{|\nabla \phi|^2} = -\frac{1}{|\nabla \phi|} (\mathbf{I} - \mathbf{n} \otimes \mathbf{n})$$

and, by symmetry of $\boldsymbol{\sigma}$, $h^1 h^2 ((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}) \cdot (\mathbf{n} \otimes \mathbf{n}) = h^1 h^2 (\mathbf{a} \cdot (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n}) \mathbf{n} = \mathbf{0}$, one thus finally obtains

$$\frac{\partial f_{el}}{\partial \nabla \phi} = -\frac{1}{|\nabla \phi|} h^1 h^2 ((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}). \quad (7.100)$$

Remark 129. The appearance of the prefactor $\frac{1}{|\nabla \phi|}$ in the derivative (7.100) is clearly somewhat problematic in regions where $|\nabla \phi|$ approaches 0 (or even worse, actually is zero). For “realistic” phase-field profiles in the obstacle case, this essentially concerns the transition region from the interface to the neighboring bulk values. Assuming a phasefield profile corresponding roughly to the standard sinusoidal shape, the degree to which this contribution can cause numerical difficulties in these regions depends quite strongly on the interpolation function used. When using any of the higher-order interpolation functions, one does not expect any actual practical difficulties since the product $h^1 h^2$ in this case converges superlinearly to zero (one factor converging to one, the other one superlinearly to zero) as one approaches the outer interface regions. In contrast, $|\nabla \phi|$ does so only linearly and one therefore expects an “extension by zero” to be a legitimate choice.

For the simplest interpolation function $h_0(\phi) = \phi$, the situation is somewhat more difficult though as $\frac{h^1 h^2}{|\nabla \phi|}$ would, for the basic one-dimensional profile, converge to a finite but non-zero value. In the discrete setting, a simple extension by zero⁶⁶ to the region just outside the interface can thus lead, in combination with the discrete divergence-operator, to a discrete Dirac-type source-term. \diamond

7.2.3 Comparison of the Different Formulations in the Two-Phase Case

Before discussing some issues related with the (inherently difficult) extension of the jump-condition based model(s) to the multiphase case, the following paragraphs will proceed to a

⁶⁵Here use is made of $\boldsymbol{\sigma}^\alpha : (\mathbf{a} \otimes \mathbf{n})_S = (\boldsymbol{\sigma}^\alpha \cdot \mathbf{a}) \cdot \mathbf{n} = (\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}) \cdot \mathbf{a}$ in order to avoid the (purely notational) difficulty of differentiating the symmetrized form with respect to a vector.

⁶⁶This may in particular happen “automatically” depending on the discretization. Using e.g. a standard cell-centered finite difference discretization based on a central gradient operator, the first cell with $\phi = 0$ would satisfy $h^1 h^2 = 0$ even though the discrete gradient would not yet vanish due to its broader stencil. This should not be interpreted as there not being a problem anymore, but simply shows that the non-zero limit $\frac{h^1 h^2}{|\nabla \phi|}$ is not properly approximated in this case.

detailed comparison of the formulations in [51], [24], [64] and [74]. Even though this is a somewhat lengthy endeavor as some of them a priori differ quite substantially in their formulation, this effort seems worthwhile for several reasons.

Firstly, this will serve to verify - as already indicated above - that the models proposed in [51] (in its small-deformation analogue), [64] (except for a missing contribution to the phasefield equation) and [74] are in fact one and the same model and just differ in their representation. While this is intuitively to be expected, this is in some cases not at all obvious based on the actual expression for the various quantities such as the elastic free energy density, the effective stress and the contributions to the phasefield equation. All three in turn can be considered to improve upon the one in [24], which, while in the same spirit, is in parts based on a somewhat different assumption.

Secondly, even though discussing different representations of an equivalent model may seem like a purely superficial effort, their differences can in fact be quite important. On the one hand, different descriptions will give rise to different implementations which can differ quite significantly in their respective computation effort. On the other hand, the various representations are suggestive of different **non-equivalent** ways of generalizing the two-phase models to the multiphase case. While some of the differences in the resulting models are fairly obvious, others are much harder to understand without properly understanding the differences in the two-phase case.

Last but not least, as each multiphase model generalization has its advantages and disadvantages, it may be quite useful to dispose of an implementation for several ones (such as is the case in the **Pace3D**-framework). As will be seen in the following discussion, the models considered do, despite their inherent complexity and the differences in the precise formulation, for the most part rely on a fairly reduced set of “central” operations. From a practical point of view, understanding which these are and making them accessible in a unified manner for the different model implementations can significantly reduce code overhead, increase readability and reduce the risk of implementation errors.

The Jump-Vector Based Approach by Mosler et al. [51]

The model in [51] is, just as the one outlined in Section 7.2.2, also based upon the use of a jump-vector \mathbf{a} and the latter can in fact be considered to be a simple consequence of the small deformation analogue of the former. Even though the details with respect to the mechanical model and the drivingforce are already worked out in Section 7.2.2, it is, in particular with respect to Remark 127, worthwhile to recall the slightly different starting point in [51] based on a variational argument instead of the (from a mechanical point of view) slightly more direct approach above.

Their approach, termed **partial rank one relaxation**, is based upon introducing a jump vector \mathbf{a} such that $[[\mathbf{F}]] = \mathbf{F}^{(2)} - \mathbf{F}^{(1)} = \mathbf{a} \otimes \mathbf{N}$ and, consistent with the averaging condition, $\mathbf{F}^1 = \mathbf{F} - \phi \mathbf{a} \otimes \mathbf{N}$ and $\mathbf{F}^2 = \mathbf{F} + (1 - \phi) \mathbf{a} \otimes \mathbf{N}$. The jump vector \mathbf{a} is then chosen such that it minimizes the total free energy for a given deformation gradient \mathbf{F} .

Translating this to the small deformation setting (i.e. in particular identifying the normal vector \mathbf{N} in the reference configuration with the one in the current one, $\mathbf{N} \approx \mathbf{n}$) and replacing the weighting by ϕ with the one in terms of the weighting-functions $h^1(\phi)$ and $h^2(\phi)$, this corresponds to defining \mathbf{F}^α and thus $\boldsymbol{\epsilon}^\alpha$, $\alpha = 1, 2$, as in Equation (7.87), where \mathbf{a} is defined in terms of the minimization problem

$$\min_{\mathbf{a}} \{f_{el}(\phi, \nabla \phi, \boldsymbol{\epsilon})\} = \min_{\mathbf{a}} \left\{ \sum_{\alpha=1}^2 h^\alpha(\phi) \frac{1}{2} (\boldsymbol{\epsilon}^\alpha(\phi, \mathbf{a}, \mathbf{n}) - \tilde{\boldsymbol{\epsilon}}^\alpha) : \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha(\phi, \mathbf{a}, \mathbf{n}) - \tilde{\boldsymbol{\epsilon}}^\alpha) \right\}. \quad (7.101)$$

Using $\boldsymbol{\epsilon}^1 = \boldsymbol{\epsilon} - h^2(\mathbf{a} \otimes \mathbf{n})_S$ and $\boldsymbol{\epsilon}^2 = \boldsymbol{\epsilon} + h^1(\mathbf{a} \otimes \mathbf{n})_S$ and differentiating with respect to \mathbf{a} leads to the FONC

$$\begin{aligned} \frac{\partial f_{el}}{\partial \mathbf{a}}(\boldsymbol{\epsilon}, \phi, \mathbf{a}, \mathbf{n}) \cdot d\mathbf{a} &= \sum_{\alpha=1}^2 h^\alpha(\phi) (\mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha(\phi, \mathbf{a}, \mathbf{n}) - \tilde{\boldsymbol{\epsilon}}^\alpha)) : \left(\frac{\partial \boldsymbol{\epsilon}^\alpha}{\partial \mathbf{a}} \cdot d\mathbf{a} \right) = \sum_{\alpha=1}^2 h^\alpha(\phi) \boldsymbol{\sigma}^\alpha : \frac{\partial \boldsymbol{\epsilon}^\alpha}{\partial \mathbf{a}} \\ &= h^1(\phi) \boldsymbol{\sigma}^1 : (-h^2(d\mathbf{a} \otimes \mathbf{n})_S) + h^2(\phi) \boldsymbol{\sigma}^2 : (h^1(d\mathbf{a} \otimes \mathbf{n})_S) \\ &= h^1(\phi) h^2(\phi) (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) : (d\mathbf{a} \otimes \mathbf{n})_S = h^1(\phi) h^2(\phi) ((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n}) \cdot d\mathbf{a} \stackrel{!}{=} \mathbf{0}, \end{aligned}$$

and thus, within the interface region (where $h^1 h^2 \neq 0$), to the continuity condition $\boldsymbol{\sigma}^1 \cdot \mathbf{n} = \boldsymbol{\sigma}^2 \cdot \mathbf{n}$. From this, it is on the one hand clear that the variational characterization by [51] is (at least whenever $h^1(\phi) h^2(\phi) \neq 0$) equivalent to imposing the jump conditions on both the strains and stresses as above, and on the other hand that the dependence of \mathbf{a} on ϕ and $\nabla \phi$ has, as already seen above, no bearing on the derivative of f_{el} with respect to ϕ and $\nabla \phi$ as these conditions characterize precisely the points where $\frac{\partial f_{el}}{\partial \mathbf{a}} = \mathbf{0}$.

An Alternative Variational Characterization of the Model by Schneider et al. [64]

The same two jump conditions are also the ones which the model by [64] and the mechanical part of the model in [24] are based upon. In contrast to the use of the jump vector \mathbf{a} above though, the description in [64] uses a transformation matrix \mathbf{Q} constructed as $\mathbf{Q}^T = \begin{pmatrix} \mathbf{n} & \mathbf{t} & \mathbf{s} \end{pmatrix}$ from an orthonormal set $(\mathbf{n}, \mathbf{t}, \mathbf{s})$ of vectors consisting of the normal \mathbf{n} and the two tangential vectors \mathbf{t} and \mathbf{s} . \mathbf{Q} is thus a unitary matrix satisfying $\mathbf{Q}^{-1} = \mathbf{Q}^T$, and can be used to transform the original entries of the stress-tensor and (analogously for the strain tensor) in their Cartesian representation to a new orthonormal coordinate system by setting

$$\boldsymbol{\sigma}_B := \mathbf{Q} \boldsymbol{\sigma} \mathbf{Q}^T = \begin{pmatrix} \mathbf{n} \cdot \boldsymbol{\sigma} \cdot \mathbf{n} & \mathbf{n} \cdot \boldsymbol{\sigma} \cdot \mathbf{t} & \mathbf{n} \cdot \boldsymbol{\sigma} \cdot \mathbf{s} \\ \mathbf{t} \cdot \boldsymbol{\sigma} \cdot \mathbf{n} & \mathbf{t} \cdot \boldsymbol{\sigma} \cdot \mathbf{t} & \mathbf{t} \cdot \boldsymbol{\sigma} \cdot \mathbf{s} \\ \mathbf{s} \cdot \boldsymbol{\sigma} \cdot \mathbf{n} & \mathbf{s} \cdot \boldsymbol{\sigma} \cdot \mathbf{t} & \mathbf{s} \cdot \boldsymbol{\sigma} \cdot \mathbf{s} \end{pmatrix} =: \begin{pmatrix} \sigma_{nn} & \sigma_{nt} & \sigma_{ns} \\ \sigma_{tn} & \sigma_{tt} & \sigma_{ts} \\ \sigma_{sn} & \sigma_{st} & \sigma_{ss} \end{pmatrix}$$

and analogously for $\boldsymbol{\epsilon}_B := \mathbf{Q} \boldsymbol{\epsilon} \mathbf{Q}^T$.

The advantage of using this coordinate system is that the matrices representing the continuity condition on the phase-specific stresses and strains therefore have a particularly simple algebraic structure as the coordinates of the normal vector \mathbf{n} in this new coordinate system are given by $\mathbf{Q}^T \mathbf{n} = \begin{pmatrix} \mathbf{n} \cdot \mathbf{n} & \mathbf{t} \cdot \mathbf{n} & \mathbf{s} \cdot \mathbf{n} \end{pmatrix}^T = \begin{pmatrix} 1 & 0 & 0 \end{pmatrix}^T$. More precisely, as $\llbracket \boldsymbol{\epsilon} \rrbracket$ is assumed to be of the form $(\mathbf{a} \otimes \mathbf{n})_S$, premultiplying this equality by \mathbf{Q} and postmultiplying the result with \mathbf{Q}^T , it is easily seen that, in this new coordinate system, $\llbracket \boldsymbol{\epsilon}_B \rrbracket$ is of the form

$$\llbracket \boldsymbol{\epsilon}_B \rrbracket = \frac{1}{2} \begin{pmatrix} \mathbf{n}^T \\ \mathbf{t}^T \\ \mathbf{s}^T \end{pmatrix} (\mathbf{a} \otimes \mathbf{n} + \mathbf{n} \otimes \mathbf{a}) \begin{pmatrix} \mathbf{n} & \mathbf{t} & \mathbf{s} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 2\mathbf{a} \cdot \mathbf{n} & \mathbf{a} \cdot \mathbf{t} & \mathbf{a} \cdot \mathbf{s} \\ \mathbf{t} \cdot \mathbf{a} & 0 & 0 \\ \mathbf{s} \cdot \mathbf{a} & 0 & 0 \end{pmatrix},$$

and, by the symmetry of the scalar product, therefore satisfies

$$\llbracket \boldsymbol{\epsilon}_B \rrbracket = \begin{pmatrix} \llbracket \boldsymbol{\epsilon}_B \rrbracket_{nn} & \llbracket \boldsymbol{\epsilon}_B \rrbracket_{nt} & \llbracket \boldsymbol{\epsilon}_B \rrbracket_{ns} \\ \llbracket \boldsymbol{\epsilon}_B \rrbracket_{nt} & 0 & 0 \\ \llbracket \boldsymbol{\epsilon}_B \rrbracket_{ns} & 0 & 0 \end{pmatrix}. \quad (7.102)$$

Remark 130. The reverse conclusion also holds, i.e. if $\boldsymbol{\epsilon}_B = \mathbf{Q} \boldsymbol{\epsilon} \mathbf{Q}^T$ is of the form as in Equation (7.102), there is some \mathbf{a} such that $\llbracket \boldsymbol{\epsilon} \rrbracket$ is of the form $(\mathbf{a} \otimes \mathbf{n})_S$. In fact, due to the symmetry of $\boldsymbol{\epsilon}_B$, obtaining such an equality requires satisfying the three equations $\mathbf{a} \cdot \mathbf{n} = \llbracket \boldsymbol{\epsilon}_{nn} \rrbracket$, $\mathbf{a} \cdot \mathbf{t} = 2\llbracket \boldsymbol{\epsilon}_{nt} \rrbracket$

and $\mathbf{a} \cdot \mathbf{s} = 2\llbracket \epsilon_{ns} \rrbracket$, or, in matrix-vector form, $\mathbf{Q}\mathbf{a} = \left(\llbracket \epsilon_{nn} \rrbracket \quad 2\llbracket \epsilon_{nt} \rrbracket \quad 2\llbracket \epsilon_{ns} \rrbracket \right)^T$. As \mathbf{Q} is a unitary matrix (with $\mathbf{Q}^{-1} = \mathbf{Q}^T$), the vector

$$\mathbf{a} = \mathbf{Q}^T \begin{pmatrix} \llbracket \epsilon_{nn} \rrbracket \\ 2\llbracket \epsilon_{nt} \rrbracket \\ 2\llbracket \epsilon_{ns} \rrbracket \end{pmatrix} \quad (7.103)$$

clearly satisfies this requirement. \diamond

Similarly, pre- and postmultiplying $\boldsymbol{\sigma}$ with \mathbf{Q} and \mathbf{Q}^T and making use of $\mathbf{n}_{\mathcal{B}} = \begin{pmatrix} 1 & 0 & 0 \end{pmatrix}^T$ and the symmetry of $\boldsymbol{\sigma}_{\mathcal{B}}$, it is easy to see that the continuity condition on the normal stresses implies that the stress jump is on the complementary entries and therefore of the form

$$\llbracket \boldsymbol{\epsilon}_{\mathcal{B}} \rrbracket = \begin{pmatrix} 0 & 0 & 0 \\ 0 & \llbracket \boldsymbol{\sigma}_{\mathcal{B}} \rrbracket_{tt} & \llbracket \boldsymbol{\sigma}_{\mathcal{B}} \rrbracket_{ts} \\ 0 & \llbracket \boldsymbol{\sigma}_{\mathcal{B}} \rrbracket_{st} & \llbracket \boldsymbol{\sigma}_{\mathcal{B}} \rrbracket_{ss} \end{pmatrix} \quad (7.104)$$

whereas the continuous components of $\boldsymbol{\sigma}$ are given by $\boldsymbol{\sigma}_{nn}$, $\boldsymbol{\sigma}_{tn} = \boldsymbol{\sigma}_{nt}$ and $\boldsymbol{\sigma}_{sn} = \boldsymbol{\sigma}_{ns}$.

Remark 131. The following discussion (in particular the retransformation to the original coordinate system in the next section) will, for shortening the notation, only be performed in a two-dimensional setting, in which, with the single tangential vector \mathbf{t} the jump in $\boldsymbol{\epsilon}$ takes the form

$$\llbracket \boldsymbol{\epsilon}_{\mathcal{B}} \rrbracket = \begin{pmatrix} \llbracket \boldsymbol{\epsilon}_{\mathcal{B}} \rrbracket_{nn} & \llbracket \boldsymbol{\epsilon}_{\mathcal{B}} \rrbracket_{nt} \\ \llbracket \boldsymbol{\epsilon}_{\mathcal{B}} \rrbracket_{nt} & 0 \end{pmatrix}$$

whereas the continuity condition on $\boldsymbol{\sigma}_{\mathcal{B}}$ carries on $\boldsymbol{\sigma}_{nn}$ and $\boldsymbol{\sigma}_{nt} = \boldsymbol{\sigma}_{tn}$. Most of it is in fact independent of whether one works in two or three dimensions though, and the rest easily carries over to the three-dimensional case as indicated in Remark 146, simply requiring somewhat more lengthy formulae. \diamond

Remark 132. The following discussion will (as in [64]) primarily make use of the in practice more common Voigt notation. Since this formalism is for the most part based upon matrix-vector and matrix-matrix multiplications only and there in particular being no point in “visually” differentiating between single- and double-contractions, single contractions will not be marked by a (\cdot) -symbol except for actual scalar products between vectors or second-order tensors in Voigt-notation.

In contrast, the contraction-symbols will continue to be used in all expressions which are to be interpreted as using the “non-Voigt”-notation. In particular, several formulae contain both Voigt-notation and the normal traction vector $\boldsymbol{\sigma} \cdot \mathbf{n}$, a notation which seems much more natural than its more cumbersome⁶⁷ counterpart $\mathbf{B}\boldsymbol{\sigma}^v$ in Voigt-notation (see below). \diamond

Using Voigt-notation, the transformations above can be reexpressed as (and similarly for the eigenstrains)

$$\boldsymbol{\sigma}_{\mathcal{B}}^v = \begin{pmatrix} \sigma_{nn} \\ \sigma_{tt} \\ \sigma_{nt} \end{pmatrix} = \mathbf{M}_{\boldsymbol{\sigma}} \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\epsilon}_{\mathcal{B}}^v = \begin{pmatrix} \epsilon_{nn} \\ \epsilon_{tt} \\ 2\epsilon_{nt} \end{pmatrix} = \mathbf{M}_{\boldsymbol{\epsilon}} \begin{pmatrix} \epsilon_{xx} \\ \epsilon_{yy} \\ 2\epsilon_{xy} \end{pmatrix}$$

using the two matrices

$$\mathbf{M}_{\boldsymbol{\sigma}} = \begin{pmatrix} n_x^2 & n_y^2 & 2n_x n_y \\ t_x^2 & t_y^2 & 2t_x t_y \\ n_x t_x & n_y t_y & n_x t_y + n_y t_x \end{pmatrix} \quad \text{and} \quad \mathbf{M}_{\boldsymbol{\epsilon}} = \begin{pmatrix} n_x^2 & n_y^2 & n_x n_y \\ t_x^2 & t_y^2 & t_x t_y \\ 2n_x t_x & 2n_y t_y & n_x t_y + n_y t_x \end{pmatrix}.$$

⁶⁷Note that the multiplication of the stress-tensor by \mathbf{n} in Voigt-notation has to map the “vector” $\boldsymbol{\sigma}^v$ to the normal traction vector, and thus (by linearity) is represented in terms of a matrix. Even though mapping the second-order tensor $\boldsymbol{\sigma}$ to the same vector in fact involves an even higher-dimensional construct, its representation in terms of the contraction with \mathbf{n} is much more compact and readable than the one in the Voigt-case.

Using this new coordinate system, $(\boldsymbol{\epsilon}^v)^\alpha = \mathbf{M}_\epsilon^{-1} \boldsymbol{\epsilon}_B^v$, and, with $\mathbf{M}_\epsilon^{-1} = \mathbf{M}_\sigma^T$, this allows rewriting the elastic free energy density of each phase as

$$\begin{aligned} f_{el}^\alpha((\boldsymbol{\epsilon}_B^v)^\alpha) &= \frac{1}{2} \left(\mathbf{M}_\sigma^T((\boldsymbol{\epsilon}_B^v)^\alpha) - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha \right) \cdot (\mathbf{C}^v)^\alpha \left(\mathbf{M}_\sigma^T((\boldsymbol{\epsilon}_B^v)^\alpha) - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha \right) \\ &= \frac{1}{2} \left((\boldsymbol{\epsilon}_B^v)^\alpha - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha \right) \cdot \left(\mathbf{M}_\sigma (\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T \right) \left((\boldsymbol{\epsilon}_B^v)^\alpha - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha \right). \end{aligned}$$

Imposing in accordance with the Hadamard condition the equality $\epsilon_{tt}^\alpha = \epsilon_{tt}$, $\alpha = 1, 2$ for the tt -component of the given strain and decomposing $(\boldsymbol{\epsilon}_B^v)^\alpha - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha$ as⁶⁸

$$(\boldsymbol{\epsilon}_B^v)^\alpha - (\tilde{\boldsymbol{\epsilon}}_B^v)^\alpha = \begin{pmatrix} \epsilon_{nn}^\alpha - \tilde{\epsilon}_{nn}^\alpha \\ \epsilon_{tt} - \tilde{\epsilon}_{tt} \\ 2(\epsilon_{nt}^\alpha - \tilde{\epsilon}_{nt}^\alpha) \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \left(\underbrace{\begin{pmatrix} \epsilon_{nn}^\alpha \\ 2\epsilon_{nt}^\alpha \end{pmatrix}}_{=:\boldsymbol{\epsilon}_n^\alpha} - \underbrace{\begin{pmatrix} \tilde{\epsilon}_{nn}^\alpha \\ 2\tilde{\epsilon}_{nt}^\alpha \end{pmatrix}}_{=:\tilde{\boldsymbol{\epsilon}}_n^\alpha} \right) + \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \left(\underbrace{\epsilon_{tt}}_{=:\epsilon_t} - \underbrace{\tilde{\epsilon}_{tt}}_{=:\tilde{\epsilon}_t} \right)$$

and using the abbreviating $\mathbf{C}_B^\alpha := (\mathbf{M}_\sigma (\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T)$ together with the symmetry of \mathbf{C}_B^α , this density can be rewritten, similarly to [64], as

$$\begin{aligned} f_{el}^\alpha((\boldsymbol{\epsilon}_B^v)^\alpha) &= \frac{1}{2} \begin{pmatrix} \epsilon_{nn}^\alpha - \tilde{\epsilon}_{nn}^\alpha \\ \epsilon_{tt} - \tilde{\epsilon}_{tt} \\ 2(\epsilon_{nt}^\alpha - \tilde{\epsilon}_{nt}^\alpha) \end{pmatrix} \cdot \mathbf{C}_B^\alpha \begin{pmatrix} \epsilon_{nn}^\alpha - \tilde{\epsilon}_{nn}^\alpha \\ \epsilon_{tt} - \tilde{\epsilon}_{tt} \\ 2(\epsilon_{nt}^\alpha - \tilde{\epsilon}_{nt}^\alpha) \end{pmatrix} \\ &= \frac{1}{2} \left((\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) \cdot \mathbf{C}_{nn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + 2(\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) \cdot \mathbf{C}_{nt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) + (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) \cdot \mathbf{C}_{tt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) \right), \end{aligned}$$

with

$$\begin{aligned} \mathbf{C}_{nn}^\alpha &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{C}_B^\alpha \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{C}_{nt}^\alpha = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{C}_B^\alpha \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \mathbf{C}_{tn}^\alpha &= \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{C}_B^\alpha \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{C}_{tt}^\alpha = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{C}_B^\alpha \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}. \end{aligned} \tag{7.105}$$

Remark 133. Note that one could also reorder the stress- and strain-entries as in [64] in order to have these auxiliary matrices correspond to contiguous subblocks of \mathbf{C}_B^α . While this does not affect the theoretical discussion, it is clearly advantageous for a practical implementation of the scheme above as this simplifies the access to the relevant part of the stiffness tensor and can be performed essentially free of overhead by appropriately permuting the rows of \mathbf{M}_ϵ and \mathbf{M}_σ during their construction.

As will be discussed below though, the primary advantage of the description in [64] (and in a somewhat different form in [24]) is that they lead (for the most part⁶⁹) to a pleasant theoretical representation. In practice, choosing such a representation in the particular coordinate system described through \mathbf{Q} above leads to a significant amount of computational overhead due to the basis transformations involved such that for actual computations, other descriptions may be preferable. \diamond

⁶⁸Even though the definition of $\boldsymbol{\epsilon}_t$ below is clearly superfluous in the two-dimensional setting considered here, it leads to a notation which is on the one hand more consistent with the one in [64] and on the other hand leads to calculations which are equally valid in the three-dimensional case where both $\boldsymbol{\epsilon}_n^\alpha$ and $\boldsymbol{\epsilon}_t$ consist of three entries.

⁶⁹An important exception here is performing the derivative with respect to $\nabla\phi$ as this, in contrast to the descriptions in [51] and [74], involves the normal and tangential vectors. In addition, a large part of these dependencies is essentially ‘‘artificial’’ and would, through a likely lengthy calculation, cancel, as their only purpose is to obtain the component representation of the vectors and matrices in this particular coordinate system, which does not affect these quantities themselves.

Fixing ϵ_t in terms of the given strain still leaves open the definition of the ϵ_n^α in terms of the given average normal strain components ϵ_n . As a first condition, one can impose that ϵ_n should be given by the weighted average of the phase-specific normal strains,

$$\epsilon_n = h^1(\phi)\epsilon_n^1 + h^2(\phi)\epsilon_n^2. \quad (7.106)$$

As a second condition, one can either, as in [64], directly impose the continuity condition $\sigma_n^1 = \sigma_n^2 = \sigma_n$, or, in a manner which underlines some of the similarities with the chemical model from Section 7.1, instead rely on an energetic argument in combination with the constraint (7.106). In fact, augmenting this constraint in the Lagrange function

$$L(\epsilon_n^1, \epsilon_n^2, \epsilon_t, \sigma_n) = f^1(\epsilon_n^1, \epsilon_n, \epsilon_t)h^1(\phi) + f^2(\epsilon_n^2, \epsilon_t)h^2(\phi) - \sigma_n \cdot (h^1(\phi)\epsilon_n^1 + h^2(\phi)\epsilon_n^2 - \epsilon_n) \quad (7.107)$$

and using the symmetry of \mathcal{C}^α , it is easily seen that the minimizers ϵ_n^1 and ϵ_n^2 for the problem

$$\min_{(\epsilon_n^1, \epsilon_n^2) : h^1(\phi)\epsilon_n^1 + h^2(\phi)\epsilon_n^2 = \epsilon_n} \left\{ f^1(\epsilon_n^1, \epsilon_t)h^1(\phi) + f^2(\epsilon_n^2, \epsilon_t)h^2(\phi) \right\}$$

are characterized by

$$\frac{\partial L}{\partial \epsilon_n^\alpha} = h^\alpha(\phi) \left(\mathcal{C}_{nn}^\alpha (\epsilon_n^\alpha - \tilde{\epsilon}_n^\alpha) + \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha) - \sigma_n \right) = \mathbf{0}, \quad \alpha = 1, 2,$$

and thus that, whenever $h^\alpha(\phi) > 0$, the decomposition of the strains with the minimal energy is characterized by the equality of the phase-specific normal stresses $\sigma_n^\alpha = \begin{pmatrix} \sigma_{nn}^\alpha & \sigma_{nt}^\alpha \end{pmatrix}^T$.

As in [64], and with the abbreviation⁷⁰ $\mathcal{S}_{nn}^\alpha = (\mathcal{C}_{nn}^\alpha)^{-1}$, one can then first solve for $(\epsilon_n^1, \epsilon_n^2)$ in terms of the (yet unknown) σ_n as $\epsilon_n^\alpha = \tilde{\epsilon}_n^\alpha + \mathcal{S}_{nn}^\alpha (\sigma_n - \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha))$ and then determine σ_n from the constraint (7.106) through

$$\sum_{\alpha=1}^2 \epsilon_n^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^2 \left(\tilde{\epsilon}_n^\alpha + \mathcal{S}_{nn}^\alpha (\sigma_n - \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha)) \right) h^\alpha(\phi) \stackrel{!}{=} \epsilon_n$$

to be given by

$$\sigma_n = \left(\sum_{\alpha=1}^2 \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right)^{-1} \left(\epsilon_n - \tilde{\epsilon}_n(\phi) + \sum_{\alpha=1}^2 \mathcal{S}_{nn}^\alpha \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha) h^\alpha(\phi) \right) \quad (7.108)$$

with $\tilde{\epsilon}_n(\phi) := \sum_{\alpha=1}^2 \tilde{\epsilon}_n^\alpha h^\alpha(\phi)$.

Reinserting this normal stress into the expressions for the ϵ_n^α and combining these with the material law then shows that the tangential stress is, as the weighted average of the σ_t^α , given by

$$\begin{aligned} \sigma_t &= \sum_{\alpha=1}^2 \left(\mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha (\sigma_n - \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha)) + \mathcal{C}_{tt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha) \right) h^\alpha(\phi) \\ &= \left(\sum_{\alpha=1}^2 \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right) \sigma_n + \sum_{\alpha=1}^2 (\mathcal{C}_{tt}^\alpha - \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha \mathcal{C}_{nt}^\alpha) (\epsilon_t - \tilde{\epsilon}_t^\alpha) h^\alpha(\phi) \\ &= \left(\sum_{\alpha=1}^2 \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right) \left(\sum_{\alpha=1}^2 \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right)^{-1} (\epsilon_n - \tilde{\epsilon}_n(\phi)) \\ &\quad + \sum_{\alpha=1}^2 (\mathcal{C}_{tt}^\alpha - \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha \mathcal{C}_{nt}^\alpha) (\epsilon_t - \tilde{\epsilon}_t^\alpha) h^\alpha(\phi) \\ &\quad + \left(\sum_{\alpha=1}^2 \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right) \left(\sum_{\alpha=1}^2 \mathcal{S}_{nn}^\alpha h^\alpha(\phi) \right)^{-1} \sum_{\alpha=1}^2 \mathcal{S}_{nn}^\alpha \mathcal{C}_{nt}^\alpha (\epsilon_t - \tilde{\epsilon}_t^\alpha) h^\alpha(\phi). \end{aligned} \quad (7.109)$$

⁷⁰Even though this notation is somewhat inconsistent as \mathcal{C}_{nn}^α is in fact the $n-n$ -subblock of \mathcal{C}^α , whereas \mathcal{S}_{nn}^α is not the $n-n$ -subblock of \mathcal{S}^α , it will nevertheless be used here as it significantly shortens the notation and there is no actual risk of confusing the meaning here as \mathcal{S}^α itself is never used.

Using the same abbreviations

$$\begin{aligned}\bar{\mathcal{T}}_{nn} &= -\sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} h^{\alpha}(\phi), & \bar{\mathcal{T}}_{nt} &= \sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nt}^{\alpha} h^{\alpha}(\phi), \\ \bar{\mathcal{T}}_{tn} &= \sum_{\alpha=1}^2 \mathcal{C}_{tn}^{\alpha} \mathcal{S}_{nn}^{\alpha} h^{\alpha}(\phi), & \bar{\mathcal{T}}_{tt} &= \sum_{\alpha=1}^2 \left(\mathcal{C}_{tt}^{\alpha} - \mathcal{C}_{tn}^{\alpha} \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nt}^{\alpha} \right) h^{\alpha}(\phi), \\ \tilde{\chi}_n &= \sum_{\alpha=1}^2 \left(\tilde{\epsilon}_n^{\alpha} + \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nt}^{\alpha} \tilde{\epsilon}_t^{\alpha} \right) h^{\alpha}(\phi), & \tilde{\chi}_t &= \sum_{\alpha=1}^2 \left(\mathcal{C}_{tt}^{\alpha} - \mathcal{C}_{tn}^{\alpha} \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nt}^{\alpha} \right) \tilde{\epsilon}_t^{\alpha}\end{aligned}\quad (7.110)$$

as in [64], it is easy to see that these stresses correspond precisely to the expressions

$$\boldsymbol{\sigma}_n = -\bar{\mathcal{T}}_{nn}^{-1} (\boldsymbol{\epsilon}_n + \bar{\mathcal{T}}_{nt} \boldsymbol{\epsilon}_t - \tilde{\chi}_n) \quad (7.111)$$

and

$$\boldsymbol{\sigma}_t = \bar{\mathcal{T}}_{tn} \boldsymbol{\sigma}_n + \bar{\mathcal{T}}_{tt} \boldsymbol{\epsilon}_t - \tilde{\chi}_t = -\bar{\mathcal{T}}_{tn} \bar{\mathcal{T}}_{nn}^{-1} (\boldsymbol{\epsilon}_n + \bar{\mathcal{T}}_{nt} \boldsymbol{\epsilon}_t - \tilde{\chi}_n) + \bar{\mathcal{T}}_{tt} \boldsymbol{\epsilon}_t - \tilde{\chi}_t \quad (7.112)$$

given there. In the form above, the calculation of the effective stress therefore consists in the following four-step procedure:

1. transform the given strain $\boldsymbol{\epsilon}^v$ and eigenstrains $(\tilde{\epsilon}^v)^{\alpha}$ to the \mathcal{B} -coordinate system;
2. transform the stiffness-tensor \mathcal{C} to the \mathcal{B} -coordinate system and calculate the submatrices in Equation (7.110);
3. calculate the stresses $\boldsymbol{\sigma}_{\mathcal{B}}$ as in equations (7.111) and (7.112);
4. retransform $\boldsymbol{\sigma}_{\mathcal{B}}^v$ to the standard coordinate system.

Retransformation of the Mechanical Model to the Original Coordinate System

Even though this is a theoretically pleasing description, it involves (as indicated above) a quite notable overhead due to the heavy use of the components of the various tensors in the transformed coordinate system. For actual computational purposes, it is thus useful to reexpress this model based on quantities in the original coordinate system only. Besides the computational point of view, this retransformation is also instructive in its own right, as it relates various operations expressed in term of \mathbf{a} - as well as for the model in [74] to be discussed below - with the more transparent “direct” calculations in the aligned coordinate-system \mathcal{B} .

For this it is more convenient to take the expression for $\boldsymbol{\sigma}_n$ in Equation (7.108) instead of Equation (7.111) as a starting point. Artificially expanding $\boldsymbol{\epsilon}_n - \tilde{\epsilon}_n(\phi)$ as $\sum_{\alpha=1}^2 (\boldsymbol{\epsilon}_n - \tilde{\epsilon}_n^{\alpha}) h^{\alpha}(\phi) = \sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nn}^{\alpha} (\boldsymbol{\epsilon}_n - \tilde{\epsilon}_n^{\alpha}) h^{\alpha}(\phi)$, $\boldsymbol{\sigma}_n$ can firstly be rewritten as

$$\boldsymbol{\sigma}_n = \left(\sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} h^{\alpha}(\phi) \right)^{-1} \sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} \left(\mathcal{C}_{nn}^{\alpha} (\boldsymbol{\epsilon}_n - \tilde{\epsilon}_n^{\alpha}) + \mathcal{C}_{nt}^{\alpha} (\boldsymbol{\epsilon}_t - \tilde{\epsilon}_t^{\alpha}) \right) h^{\alpha}(\phi), \quad (7.113)$$

where the term in the latter parenthesis is precisely the normal phase-specific stress component corresponding to the Voigt-Taylor model, i.e. denoting these stresses by

$$\boldsymbol{\Sigma}_n^{\alpha} := \mathcal{C}_{nn}^{\alpha} (\boldsymbol{\epsilon}_n - \tilde{\epsilon}_n^{\alpha}) + \mathcal{C}_{nt}^{\alpha} (\boldsymbol{\epsilon}_t - \tilde{\epsilon}_t^{\alpha}),$$

one has

$$\boldsymbol{\sigma}_n = \left(\sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} h^{\alpha}(\phi) \right)^{-1} \sum_{\alpha=1}^2 \mathcal{S}_{nn}^{\alpha} \boldsymbol{\Sigma}_n^{\alpha} h^{\alpha}(\phi). \quad (7.114)$$

Based upon the middle expression for $\boldsymbol{\sigma}_t$,

$$\boldsymbol{\sigma}_t = \left(\sum_{\alpha=1}^2 \mathcal{C}_{tn}^{\alpha} \mathcal{S}_{nn}^{\alpha} h^{\alpha}(\phi) \right) \boldsymbol{\sigma}_n + \sum_{\alpha=1}^2 \left(\mathcal{C}_{tt}^{\alpha} - \mathcal{C}_{tn}^{\alpha} \mathcal{S}_{nn}^{\alpha} \mathcal{C}_{nt}^{\alpha} \right) (\boldsymbol{\epsilon}_t - \tilde{\epsilon}_t^{\alpha}) h^{\alpha}(\phi)$$

in Equation (7.109) (or from the first one in (7.112) and the definitions in Equation (7.110)), one can derive an expression of similar nature for the tangential stress components. Adding and subtracting $\sum_{\alpha=1}^2 \mathbf{C}_{tn}^{\alpha}(\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha})$, the second term (corresponding to $\bar{\mathbf{T}}_{tt}\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\chi}}_t$ in the notation of [64]) can be expanded to

$$\begin{aligned} & \sum_{\alpha=1}^2 \left(\mathbf{C}_{tt}^{\alpha}(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^{\alpha}) + \mathbf{C}_{tn}^{\alpha}(\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha}) \right) h^{\alpha}(\phi) - \sum_{\alpha=1}^2 \mathbf{C}_{tn}^{\alpha} \left((\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha}) + \mathbf{S}_{nn}^{\alpha} \mathbf{C}_{nt}^{\alpha}(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^{\alpha}) \right) h^{\alpha}(\phi) \\ &= \sum_{\alpha=1}^2 \left(\mathbf{C}_{tt}^{\alpha}(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^{\alpha}) + \mathbf{C}_{tn}^{\alpha}(\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha}) \right) h^{\alpha}(\phi) - \sum_{\alpha=1}^2 \mathbf{C}_{tn}^{\alpha} \mathbf{S}_{nn}^{\alpha} \left(\mathbf{C}_{nn}^{\alpha}(\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha}) + \mathbf{C}_{nt}^{\alpha}(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^{\alpha}) \right) h^{\alpha}(\phi). \end{aligned}$$

The term in the last parenthesis is just $\boldsymbol{\Sigma}_n^{\alpha}$, whereas the first one is nothing but the tangential component of the stress predicted by the Voigt-Taylor model, i.e. with the analogous abbreviation $\boldsymbol{\Sigma}_t^{\alpha} := \mathbf{C}_{tt}^{\alpha}(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^{\alpha}) + \mathbf{C}_{tn}^{\alpha}(\boldsymbol{\epsilon}_n^{\alpha} - \tilde{\boldsymbol{\epsilon}}_n^{\alpha})$ as for the normal components, one has

$$\bar{\mathbf{T}}_{tt}\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\chi}}_t = \sum_{\alpha=1}^2 \left(\boldsymbol{\Sigma}_t^{\alpha} - \mathbf{C}_{tn}^{\alpha} \mathbf{S}_{nn}^{\alpha} \boldsymbol{\Sigma}_n^{\alpha} \right) h^{\alpha}(\phi).$$

The prefactors $\mathbf{C}_{tn}^{\alpha} \mathbf{S}_{nn}^{\alpha}$ in the last sum are the same as the ones appearing in the prefactor of $\bar{\mathbf{T}}_{tn}\boldsymbol{\sigma}_n$, allowing to summarize the formula for the tangential stress components to

$$\boldsymbol{\sigma}_t = \sum_{\alpha=1}^2 \boldsymbol{\Sigma}_t^{\alpha} h^{\alpha}(\phi) + \sum_{\alpha=1}^2 \mathbf{C}_{tn}^{\alpha} \mathbf{S}_{nn}^{\alpha} (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^{\alpha}) h^{\alpha}(\phi), \quad (7.115)$$

i.e. a first Voigt-Taylor-type contribution corrected by a term based on the differences of the real normal stress from the one predicted by the Voigt-Taylor model.

Remark 134. Even though this may not seem particularly convenient at this point, one can also artificially rewrite $\boldsymbol{\sigma}_n$ in a similar form. In fact, simply expanding expressions by making use of $\sum_{\alpha=1}^2 h^{\alpha}(\phi) = 1$ and $\mathbf{C}_{nn}^{\alpha} \mathbf{S}_{nn}^{\alpha} = \mathbf{I}$, one has

$$\boldsymbol{\sigma}_n = \sum_{\alpha=1}^2 \boldsymbol{\Sigma}_n^{\alpha} h^{\alpha}(\phi) + \sum_{\alpha=1}^2 (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^{\alpha}) h^{\alpha}(\phi) = \sum_{\alpha=1}^2 \boldsymbol{\Sigma}_n^{\alpha} h^{\alpha}(\phi) + \sum_{\alpha=1}^2 \mathbf{C}_{nn}^{\alpha} \mathbf{S}_{nn}^{\alpha} (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^{\alpha}) h^{\alpha}(\phi). \quad (7.116)$$

This obviously does not actually contain any additional information (and in particular does not allow for the determination of $\boldsymbol{\sigma}_n$ itself, for which one has to rely on Equation (7.114)), but shows that both the tangential and normal stress components can be expressed in a common form once the actual normal stress $\boldsymbol{\sigma}_n$ is known, namely a Voigt-Taylor-type prediction together with a correction based on the $\mathbf{S}_{nn}^{\alpha}(\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^{\alpha})$, $\alpha = 1, 2$ and \mathbf{C}_{tn}^{α} resp. \mathbf{C}_{nn}^{α} . \diamond

Despite of the relative simplicity of the formulae for $\boldsymbol{\sigma}_n$ and $\boldsymbol{\sigma}_t$ when expressed in terms of the auxiliary quantities $\boldsymbol{\Sigma}_{n/t}^{\alpha}$, both equations (7.114) and (7.115) still rely heavily upon the transformed quantities (and in particular the transformed stiffness-tensors $\mathbf{C}_{\mathcal{B}}^{\alpha}$) in the \mathcal{B} -system. A large part of this dependence can in fact be eliminated in a straightforward manner though as a closer look at the interplay between the extraction of the normal components and the transformation matrices will reveal.

A convenient starting point here are the \mathbf{C}_{nn}^{α} -matrices, which, by Equation (7.105), are determined by

$$\mathbf{C}_{nn}^{\alpha} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{C}_{\mathcal{B}}^{\alpha} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \left(\mathbf{M}_{\boldsymbol{\sigma}} (\mathbf{C}^v)^{\alpha} \mathbf{M}_{\boldsymbol{\sigma}}^T \right) \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

Making use of

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{M}_{\boldsymbol{\sigma}} = \begin{pmatrix} n_x^2 & n_y^2 & 2n_x n_y \\ n_x t_x & n_y t_y & n_x t_y + n_y t_x \end{pmatrix} = \underbrace{\begin{pmatrix} n_x & n_y \\ t_x & t_y \end{pmatrix}}_{=: \mathbf{Q}} \begin{pmatrix} n_x & 0 & n_y \\ 0 & n_y & n_x \end{pmatrix}$$

and similarly

$$M_{\sigma}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} n_x & 0 \\ 0 & n_y \\ n_y & n_x \end{pmatrix} Q^T$$

this can also be rewritten as

$$\mathbf{C}_{nn}^{\alpha} = Q \left(\begin{pmatrix} n_x & 0 & n_y \\ 0 & n_y & n_x \end{pmatrix} (\mathbf{C}^v)^{\alpha} \begin{pmatrix} n_x & 0 \\ 0 & n_y \\ n_y & n_x \end{pmatrix} \right) Q^T.$$

Adopting the same notation \mathbf{B} resp. \mathbf{B}^T as in [74] (not to be confused with the caligraphic indicator \mathcal{B} of the basis),

$$\mathbf{B} := \begin{pmatrix} n_x & 0 & n_y \\ 0 & n_y & n_x \end{pmatrix}, \quad (7.117)$$

these relations can be written more succinctly as

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} M_{\sigma} = Q \mathbf{B}, \quad M_{\sigma}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{B}^T Q^T \quad \text{and} \quad \mathbf{C}_{nn}^{\alpha} = Q \mathbf{B} (\mathbf{C}^v)^{\alpha} \mathbf{B}^T Q^T. \quad (7.118)$$

Remark 135. It is easy to verify that $\mathbf{B}(\mathbf{C}^v)^{\alpha}\mathbf{B}^T$ corresponds, in non-Voigt notation, to the matrix $\mathbf{n} \cdot \mathbf{C}^{\alpha} \cdot \mathbf{n}$. In addition, while the use of \mathbf{B} is convenient for notational purposes, this matrix need not actually be constructed in practice. Instead, it is both more “readable” and more efficient to evaluate the resulting expressions directly (i.e. the first entry of $\mathbf{C}^v \mathbf{B}^T$ simply corresponds to $C^{xxxx}n_x + C^{xxyy}n_y$ in 2D resp. $C^{xxxx}n_x + C^{xxyy}n_y + C^{xxzz}n_z$ in 3-D, which is easily translated to Voigt-notation). \diamond

Remark 136. The role of the outer factors \mathbf{Q} and \mathbf{Q}^T here is precisely the same as in [64], namely to transform the inner quantity from the Cartesian coordinate system to the system \mathcal{B} . In contrast to the original procedure of first transforming the fourth-order tensors $(\mathbf{C}^v)^{\alpha}$ in Voigt-notation using M_{σ} and $M_{\epsilon}^{-1} = M_{\sigma}^T$ and then extracting the normal subblocks, the major advantage of first contracting \mathbf{C}^{α} with \mathbf{n} from the left and right (resp. left-right-multiplying $(\mathbf{C}^v)^{\alpha}$ with \mathbf{B} and \mathbf{B}^T) is that $\mathbf{n} \cdot \mathbf{C}^{\alpha} \cdot \mathbf{n}$ is a “standard” second-order tensor, which, even though it is symmetric, is not subject to the pitfalls of the Voigt-notation. In particular, it transforms based on the standard rules using \mathbf{Q} and \mathbf{Q}^T , instead of the more complex original transformation. Being composed of the orthonormal vectors \mathbf{n} and \mathbf{t} , it is also clear that \mathbf{Q} is unitary, i.e. satisfies $\mathbf{Q}\mathbf{Q}^T = \mathbf{Q}^T\mathbf{Q} = \mathbf{I}$. \diamond

With the outer matrices in the expression for \mathbf{C}_{nn}^{α} being inverses of each other, this leads to the alternative expression

$$\mathbf{S}_{nn}^{\alpha} = \mathbf{Q} \left(\mathbf{B} (\mathbf{C}^v)^{\alpha} \mathbf{B}^T \right)^{-1} \mathbf{Q}^T \quad (7.119)$$

for \mathbf{S}_{nn}^{α} and thus $\sigma_{\mathbf{n}}$ from Equation (7.114) can, making use of $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$, be rewritten as

$$\sigma_{\mathbf{n}} = \mathbf{Q} \left(\sum_{\alpha=1}^2 \left(\mathbf{B} (\mathbf{C}^v)^{\alpha} \mathbf{B}^T \right)^{-1} h^{\alpha}(\phi) \right)^{-1} \sum_{\alpha=1}^2 \left(\mathbf{B} (\mathbf{C}^v)^{\alpha} \mathbf{B}^T \right)^{-1} \mathbf{Q}^T \Sigma_{\mathbf{n}}^{\alpha} h^{\alpha}(\phi).$$

While the prefactor itself does not now make explicit use of the transformed stiffness-tensors $\mathbf{C}_{\mathcal{B}}^{\alpha}$, the evaluation of the $\Sigma_{\mathbf{n}}^{\alpha}$ itself still does. In accordance with the physical meaning of $\Sigma_{\mathbf{n}}^{\alpha}$, it is not surprising that one is also able to reexpress $\Sigma_{\mathbf{n}}^{\alpha}$ based on simple quantities from the original coordinate system. In fact, from the original expression $\Sigma_{\mathbf{n}} = \mathbf{C}_{nn}^{\alpha}(\epsilon_{\mathbf{n}} - \tilde{\epsilon}_{\mathbf{n}}^{\alpha}) + \mathbf{C}_{nt}^{\alpha}(\epsilon_{\mathbf{t}} - \tilde{\epsilon}_{\mathbf{t}}^{\alpha})$, the

definitions of $\boldsymbol{\epsilon}_n - \tilde{\boldsymbol{\epsilon}}_n^\alpha$ and $\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t$ and again making use of Equation (7.118) together with (7.105), $\boldsymbol{\Sigma}_n^\alpha$ may be rewritten as

$$\begin{aligned} \boldsymbol{\Sigma}_n^\alpha = & \left(\mathbf{Q} \mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \right) \left(\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{M}_\epsilon(\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha) \right) \\ & + \left(\mathbf{Q} \mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right) \left(\begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\epsilon(\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha) \right). \end{aligned}$$

By extracting both the common pre- and postfactors $\mathbf{Q} \mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T$ and $\mathbf{M}_\epsilon(\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha)$ and since the sum of the two inner dyads is clearly the (3×3) identity matrix, this simplifies first to $\boldsymbol{\Sigma}_n^\alpha = \mathbf{Q} \mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{M}_\sigma^T \mathbf{I} \mathbf{M}_\epsilon(\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha)$ and then, with $\mathbf{M}_\sigma^T = \mathbf{M}_\epsilon^{-1}$, to

$$\boldsymbol{\Sigma}_n^\alpha = \mathbf{Q} \mathbf{B}(\mathbf{C}^v)^\alpha (\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha) = \mathbf{Q} \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha \quad (7.120)$$

with $(\boldsymbol{\Sigma}^v)^\alpha := (\mathbf{C}^v)^\alpha (\boldsymbol{\epsilon}^v - (\tilde{\boldsymbol{\epsilon}}^v)^\alpha)$. Reinserting this into the expression for $\boldsymbol{\sigma}_n$ and again cancelling \mathbf{Q}^T and \mathbf{Q} , one finally has

$$\boldsymbol{\sigma}_n = \mathbf{Q} \left(\sum_{\alpha=1}^2 (\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} h^\alpha(\phi) \right)^{-1} \sum_{\alpha=1}^2 (\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi), \quad (7.121)$$

i.e. an expression which, except for the appearance of the \mathbf{Q} -prefactor in front, is entirely based upon non-transformed quantities.

Remark 137. Instead of interpreting $\boldsymbol{\sigma}_n$ as the sub-vector of normal entries of the stress-tensor in Voigt-notation, one can also interpret it as the components of the normal stress $\boldsymbol{\sigma} \cdot \mathbf{n}$ in the system \mathcal{B} . In fact, as

$$\boldsymbol{\sigma}_n = \begin{pmatrix} \sigma_{nn} \\ \sigma_{tn} \end{pmatrix} = \begin{pmatrix} \mathbf{n} \cdot \boldsymbol{\sigma} \cdot \mathbf{n} \\ \mathbf{t} \cdot \boldsymbol{\sigma} \cdot \mathbf{n} \end{pmatrix} = \mathbf{Q}(\boldsymbol{\sigma} \cdot \mathbf{n}),$$

it is clear that the role of \mathbf{Q} (being composed row-wise of the two orthonormal vectors \mathbf{n} and \mathbf{t}) in Equation (7.121) is simply to “read out” the coordinates of $\boldsymbol{\sigma} \cdot \mathbf{n}$ in the basis formed by (\mathbf{n}, \mathbf{t}) . Applying \mathbf{Q}^T from the left (i.e. evaluating $\sigma_{nn} \mathbf{n} + \sigma_{tn} \mathbf{t}$), it follows that the actual normal stress is given by

$$\begin{aligned} \boldsymbol{\sigma} \cdot \mathbf{n} = & \left(\sum_{\alpha=1}^2 (\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} h^\alpha(\phi) \right)^{-1} \sum_{\alpha=1}^2 (\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi) \\ = & \left(\sum_{\alpha=1}^2 (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} h^\alpha(\phi) \right)^{-1} \cdot \sum_{\alpha=1}^2 (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \mathbf{n}) h^\alpha(\phi) \end{aligned} \quad (7.122)$$

where $\boldsymbol{\Sigma}^\alpha := \mathbf{C}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)$ and can in particular be evaluated completely without any reference to the transformed system \mathcal{B} . \diamond

Remark 138. While it is tempting, based on the role of \mathbf{Q} and \mathbf{Q}^T , to directly try to retransform Equation (7.121) back into the corresponding part of the stress in the standard coordinate system, this is not as straightforward as one might expect at first sight. In fact, all that would be required in order to do so would be to “reinsert” the entries of $\boldsymbol{\sigma}_n$ into the full vector $\boldsymbol{\sigma}_\mathcal{B}$ and then, as the inverse of \mathbf{M}_σ is given by \mathbf{M}_ϵ^T , apply \mathbf{M}_ϵ^T from the left. In matrix-vector form,

this corresponds to applying

$$\begin{aligned} \mathbf{M}_\epsilon^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} &= \begin{pmatrix} n_x^2 & t_x^2 & 2n_x t_x \\ n_y^2 & t_y^2 & 2n_y t_y \\ n_x n_y & t_x t_y & n_x t_y + n_y t_x \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} n_x^2 & 2n_x t_x \\ n_y^2 & 2n_y t_y \\ n_x n_y & n_x t_y + n_y t_x \end{pmatrix} \\ &= \begin{pmatrix} n_x & 0 \\ 0 & n_y \\ \frac{1}{2}n_y & \frac{1}{2}n_x \end{pmatrix} \begin{pmatrix} n_x & 2t_x \\ n_y & 2t_y \end{pmatrix} = \begin{pmatrix} n_x & 0 \\ 0 & n_y \\ \frac{1}{2}n_y & \frac{1}{2}n_x \end{pmatrix} \begin{pmatrix} n_x & t_x \\ n_y & t_y \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \end{aligned} \quad (7.123)$$

to the expression for $\boldsymbol{\sigma}_n$. Even though the middle matrix is just \mathbf{Q}^T , there is an additional “weighting” by the right-most matrix which prohibits cancelling \mathbf{Q}^T and the left-most factor \mathbf{Q} in Equation (7.121).

One might expect at first sight that the factors 2 in the right-most matrix and the $\frac{1}{2}$ -prefactors in the left-most matrix are due to the Voigt-representation, but this is actually not the case (note that the transformation above is between two stress-type quantities). In fact (by the previous remark) $\boldsymbol{\sigma}_n$ in Equation (7.121) corresponds to the correct components of the normal traction vector $\boldsymbol{\sigma}_n$ expressed in the \mathcal{B} -system, and retransforming this **vector** back to the Cartesian system could therefore indeed be done by applying \mathbf{Q}^T from the left. Interpreted as the components of a symmetric second-order tensor written in Voigt-notation, the transformation is somewhat different though, a point which will be discussed in a little more detail in Remarks 140 and 141 below. \diamond

Using this expression for $\boldsymbol{\sigma}_n$, one can now evaluate the tangential stress components $\boldsymbol{\sigma}_t$ in Equation (7.115). On the one hand, arguing analogously as for $\boldsymbol{\Sigma}_n^\alpha$ in Equation (7.120), it is easy to verify that $\boldsymbol{\Sigma}_t^\alpha = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma (\boldsymbol{\Sigma}^v)^\alpha$. On the other hand, combining the expression for \mathcal{C}_{tn}^α from Equation (7.105) with the definition of $\mathcal{C}_\mathcal{B}^\alpha$ and using (7.118), one has

$$\mathcal{C}_{tn}^\alpha = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma (\mathcal{C}^v)^\alpha \mathbf{B} \mathbf{Q}^T.$$

and thus with Equation (7.119)

$$\mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma (\mathcal{C}^v)^\alpha \mathbf{B} (\mathbf{B} (\mathcal{C}^v)^\alpha \mathbf{B}^T)^{-1} \mathbf{Q}^T.$$

Inserting these expressions into the formula (7.115) for $\boldsymbol{\sigma}_t$, one obtains

$$\begin{aligned} \boldsymbol{\sigma}_t &= \sum_{\alpha=1}^2 \boldsymbol{\Sigma}_t^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 \mathcal{C}_{tn}^\alpha \mathcal{S}_{nn}^\alpha (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^\alpha) h^\alpha(\phi) \\ &= \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma \left(\sum_{\alpha=1}^2 (\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 (\mathcal{C}^v)^\alpha \mathbf{B} (\mathbf{B} (\mathcal{C}^v)^\alpha \mathbf{B}^T)^{-1} \mathbf{Q}^T (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^\alpha) h^\alpha(\phi) \right). \end{aligned}$$

By the previous Remark 138, $\mathbf{Q}^T (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^\alpha)$ is nothing but $\boldsymbol{\sigma} \cdot \mathbf{n} - \boldsymbol{\Sigma}^\alpha \cdot \mathbf{n}$, i.e. an expression which can be evaluated fully without recourse to the system \mathcal{B} . As the other terms in the last factor do not rely on \mathcal{B} either, the only remaining fragment of the transformation consists in the prefactor $\begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma$, whose role is simply reading out the tt -component of the expression on the right⁷¹.

It turns out that, with respect to the construction of the actual stress $\boldsymbol{\sigma}^v$, this is now completely

⁷¹This is clear by construction, but can also be verified directly. In fact,

$$\begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma = \begin{pmatrix} t_x^2 & t_y^2 & 2t_x t_y \end{pmatrix} = \begin{pmatrix} t_x & t_y \end{pmatrix} \begin{pmatrix} t_x & 0 & t_y \\ 0 & t_y & t_x \end{pmatrix}.$$

It is, as for \mathbf{B} , easily verified that the pre-multiplication of a stress-type tensor $\hat{\boldsymbol{\sigma}}^v$ in Voigt-notation by second factor simply corresponds to the evaluation of $\hat{\boldsymbol{\sigma}} \cdot \mathbf{t}$. Further applying $\begin{pmatrix} t_x & t_y \end{pmatrix}$ then leads to $\mathbf{t} \cdot \hat{\boldsymbol{\sigma}} \cdot \mathbf{t} = \hat{\boldsymbol{\sigma}}_{tt}$.

superfluous as the right factor in fact already represents the (full) correct stress. Eliminating $\boldsymbol{\sigma}_n$ and $\boldsymbol{\Sigma}_n^\alpha$ (recall that these quantities are still based upon the system \mathcal{B}) with $\boldsymbol{\sigma}_n = \mathbf{QB}\boldsymbol{\sigma}^v$ (see Remark 137 and recall that a multiplication by \mathbf{B} in Voigt-notation corresponds to the multiplication by \mathbf{n}) and $\boldsymbol{\Sigma}_n^\alpha = \mathbf{QB}(\boldsymbol{\Sigma}^v)^\alpha$, one therefore has the formula

$$\boldsymbol{\sigma}^v = \sum_{\alpha=1}^2 (\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 (\mathbf{C}^v)^\alpha \mathbf{B}(\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} ((\boldsymbol{\sigma} \cdot \mathbf{n}) - \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha) h^\alpha(\phi), \quad (7.124)$$

resp., in non-Voigt notation,

$$\boldsymbol{\sigma} = \sum_{\alpha=1}^2 \boldsymbol{\Sigma}^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 (\mathbf{C}^\alpha \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot ((\boldsymbol{\sigma} - \boldsymbol{\Sigma}^\alpha) \cdot \mathbf{n}) h^\alpha(\phi), \quad (7.125)$$

with the normal traction vector $\boldsymbol{\sigma} \cdot \mathbf{n}$ given by Equation (7.122).

That this is true for the tangential part of the stress is clear from the above. For the normal part, this can either be seen from the ‘‘artificial’’ reexpression of $\boldsymbol{\sigma}_n$ in Equation (7.116) and inserting the expressions for \mathbf{C}_{nn}^α and \mathbf{S}_{nn}^α , or, more succinctly, by verifying that the expression (7.124) reproduces the correct remaining two (normal) components of the stress. This is easily seen to be the case, as extracting the normal componts in the \mathcal{B} -system corresponds, by Equation (7.118), to a premultiplication by \mathbf{QB} , from which one obtains

$$\mathbf{QB}\boldsymbol{\sigma}^v = \mathbf{QB} \sum_{\alpha=1}^2 (\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 \mathbf{QB}(\mathbf{C}^v)^\alpha \mathbf{B}(\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}^T)^{-1} ((\boldsymbol{\sigma} \cdot \mathbf{n}) - \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha) h^\alpha(\phi).$$

As $\mathbf{B}(\mathbf{C}^v)^\alpha \mathbf{B}$ and its inverse cancel in the second term, this leaves

$$\mathbf{QB}\boldsymbol{\sigma}^v = \sum_{\alpha=1}^2 \mathbf{QB}(\boldsymbol{\Sigma}^v)^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^2 \mathbf{Q}((\boldsymbol{\sigma} \cdot \mathbf{n}) - \mathbf{B}(\boldsymbol{\Sigma}^v)^\alpha) h^\alpha(\phi) = \sum_{\alpha=1}^2 \mathbf{Q}(\boldsymbol{\sigma} \cdot \mathbf{n}) h^\alpha(\phi) = \mathbf{Q}(\boldsymbol{\sigma} \cdot \mathbf{n}).$$

According to Remark 137, $\mathbf{Q}(\boldsymbol{\sigma} \cdot \mathbf{n}) = \boldsymbol{\sigma}_n$, from which it follows that Equation (7.124) (resp. Equation (7.125)) reproduces the correct values of all three entries of $\boldsymbol{\sigma}_\mathcal{B}^v$ in the \mathcal{B} -system, and therefore also, after a retransformation to the original system, of $\boldsymbol{\sigma}^v$.

Remark 139. Note that the expressions (7.124) and (7.125) are now entirely based upon quantities in the non-transformed system. While this is of course still the same model, it has the major advantage of avoiding the construction and repeated application of the transformation matrices \mathbf{M}_σ and \mathbf{M}_ϵ (resp. their transposes). Besides a reduction in computational cost and complexity, this in particular also avoids the construction of a (resp. two in 3D) tangential vector. This is particularly helpful when trying to evaluate $\frac{\partial f_{el}}{\partial \nabla \phi}$, as this contribution now only enters (in an explicit form) in terms of \mathbf{n} , instead of in a quite implicit form through both the normal and the tangent vector(s) in the original formulation. \diamond

Remark 140. As seen above, the expressions in equations (7.124) resp. (7.125) - a priori only necessary for the tangential part of the stress-tensor, as the normal stresses are already known - allow reconstructing the full stress tensor without an explicit separate treatment of the normal and tangential components. It is nevertheless interesting to understand the meaning of the stress-tensor (resp. vector in Voigt-notation) obtained by retransforming the expressions for $\boldsymbol{\sigma}_n$ in Equations (7.120) resp. (7.121) as indicated in Remark 138.

The central observation here is that the combined operation of reinserting the normal stress-components into $\boldsymbol{\sigma}_\mathcal{B}^v$ and then retransforming this vector can also be rewritten as

$$\boldsymbol{\sigma}_n^v = \mathbf{M}_\epsilon^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \boldsymbol{\sigma}_n = \sigma_{nn} \begin{pmatrix} n_x^2 \\ n_y^2 \\ n_x n_y \end{pmatrix} + \sigma_{nt} \begin{pmatrix} 2n_x t_x \\ 2n_y t_y \\ n_x t_y + n_y t_x \end{pmatrix},$$

corresponding to the Voigt-representation of $\sigma_{nn}\mathbf{n} \otimes \mathbf{n} + \sigma_{nt}\mathbf{n} \otimes \mathbf{t} + \sigma_{nt}\mathbf{t} \otimes \mathbf{n}$, i.e. the symmetric second-order tensor with the same normal components as $\sigma_{\mathbf{n}}$ and a vanishing tangential one. From this, it is then also clear that extracting the normal components from $\sigma_{\mathbf{n}}^v$ (in contrast to $\sigma_{\mathbf{n}}$ which is a full tensor in Voigt-notation in the standard coordinate system) and reinserting them in the same manner will again lead to the same tensor $\sigma_{\mathbf{n}}^v$. This implies that, given any symmetric second-order tress-type tensor σ^v (in Voigt-notation), the matrix

$$N_{\sigma}^v := M_{\epsilon}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} M_{\sigma} = M_{\epsilon}^T \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} M_{\sigma}$$

describing this combined operation corresponds to a projection operator onto the “normal subspace” of the stress-type tensors, while the matrix

$$T_{\sigma}^v := I - N_{\sigma}^v = M_{\epsilon}^T I M_{\sigma} - N_{\sigma}^v = M_{\epsilon}^T \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} M_{\sigma}$$

corresponds to the projection operator onto the “tangent subspace” of the stress-type tensors⁷². The stress obtained by retransforming $\sigma_{\mathbf{n}}$ as above is therefore simply the “minimal” one with the prescribed normal component $\sigma \cdot \mathbf{n}$. \diamond

Besides the calculation rule for σ^v , both the driving force calculation and dealing with more complex models - such as visco-elastic or plastic settings - will also require access to the phase-specific strains and stresses. Their calculation is again straightforward in the \mathcal{B} -system as, by assumption, the tangential ones are all equal to the average ones, $\epsilon_{\mathbf{t}}^{\alpha} = \epsilon_{\mathbf{t}}$, while the normal ones can be recovered based upon

$$\sigma_{\mathbf{n}} = \mathcal{C}_{nn}^{\alpha}(\epsilon_{\mathbf{n}}^{\alpha} - \tilde{\epsilon}_{\mathbf{n}}^{\alpha}) + \mathcal{C}_{nt}^{\alpha}(\epsilon_{\mathbf{t}} - \tilde{\epsilon}_{\mathbf{t}}^{\alpha}) = \mathcal{C}_{nn}^{\alpha}(\epsilon_{\mathbf{n}}^{\alpha} - \epsilon_{\mathbf{n}}) + \mathcal{C}_{nn}^{\alpha}(\epsilon_{\mathbf{n}} - \tilde{\epsilon}_{\mathbf{n}}^{\alpha}) + \mathcal{C}_{nt}^{\alpha}(\epsilon_{\mathbf{t}} - \tilde{\epsilon}_{\mathbf{t}}^{\alpha})$$

and thus

$$\epsilon_{\mathbf{n}}^{\alpha} - \epsilon_{\mathbf{n}} = \mathcal{S}_{nn}^{\alpha}(\sigma_{\mathbf{n}} - \mathcal{C}_{nt}^{\alpha}(\epsilon_{\mathbf{t}} - \tilde{\epsilon}_{\mathbf{t}}^{\alpha}) - \mathcal{C}_{nn}^{\alpha}(\epsilon_{\mathbf{n}} - \tilde{\epsilon}_{\mathbf{n}}^{\alpha})) = \mathcal{S}_{nn}^{\alpha}(\sigma_{\mathbf{n}} - \Sigma_{\mathbf{n}}^{\alpha}). \quad (7.126)$$

Combined with the equality of the tangential stresses, $\epsilon_{\mathbf{t}}^{\alpha} - \epsilon_{\mathbf{t}} = \mathbf{0}$, one therefore has, in the \mathcal{B} -system (again using Equation (7.118)),

$$\epsilon_{\mathcal{B}}^{\alpha} - \epsilon_{\mathcal{B}} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \mathcal{S}_{nn}^{\alpha}(\sigma_{\mathbf{n}} - \Sigma_{\mathbf{n}}^{\alpha}) = \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} Q \left(B(\mathcal{C}^v)^{\alpha} B^T \right)^{-1} \left(Q^T \sigma_{\mathbf{n}} - B(\Sigma^v)^{\alpha} \right).$$

In contrast to the recovery of the normal stresses in the original coordinate system, the recovery of this purely normal jump in the strains results in a simple formula. In fact, using $M_{\epsilon}^{-1} = M_{\sigma}^T$

⁷²One way to see this is by realizing that the inner matrices $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ and $\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$ in the definitions

of N_{σ}^v and T_{σ}^v are clearly the corresponding projectors in the \mathcal{B} -system with the M_{σ} - and M_{ϵ}^T -matrices just transforming to and back from that system.

Alternatively, it can be verified in a more explicit but also more tedious manner in the original system. A quite useful property in this respect is that, as one can verify making use of the orthogonality properties of \mathbf{n} and \mathbf{t} ,

$$B M_{\epsilon}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} = Q^T,$$

which in particular ensures that, in combination with Equation (7.118), $B N_{\sigma}^v = Q^T Q B = B$, and thus the equality of the normal stress vectors $B \sigma^v$ and $B N_{\sigma}^v \sigma^v$ before and after the projection.

and combining $M_{\sigma}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{B}^T \mathbf{Q}^T$ (see Equation (7.118)) with $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ and $\mathbf{Q}^T \boldsymbol{\sigma}_n = \boldsymbol{\sigma} \cdot \mathbf{n}$, the retransformation leads to

$$(\boldsymbol{\epsilon}^v)^\alpha - \boldsymbol{\epsilon}^v = \mathbf{B}^T \left(\mathbf{B} (\mathbf{C}^v)^\alpha \mathbf{B}^T \right)^{-1} \left((\boldsymbol{\sigma} \cdot \mathbf{n}) - \mathbf{B} (\boldsymbol{\Sigma}^v)^\alpha \right). \quad (7.127)$$

Defining

$$\mathbf{a}^\alpha := \left(\mathbf{B} (\mathbf{C}^v)^\alpha \mathbf{B}^T \right)^{-1} \left((\boldsymbol{\sigma} \cdot \mathbf{n}) - \mathbf{B} (\boldsymbol{\Sigma}^v)^\alpha \right) \quad \text{resp.} \quad \mathbf{a}^\alpha := (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\sigma} - \boldsymbol{\Sigma}^\alpha) \cdot \mathbf{n} \right) \quad (7.128)$$

this can also, in a form coinciding with the jump vector based formulation outlined above, be written as

$$(\boldsymbol{\epsilon}^v)^\alpha - \boldsymbol{\epsilon}^v = \mathbf{B}^T \mathbf{a} = \begin{pmatrix} n_x & 0 \\ 0 & n_y \\ n_y & n_x \end{pmatrix} \mathbf{a} = \begin{pmatrix} a_x n_x \\ a_y n_y \\ a_x n_y + a_y n_x \end{pmatrix}, \quad (7.129)$$

with is easily recognized as the Voigt-representation of the strain-type (i.e. with the shear components doubled) tensor $\boldsymbol{\epsilon}^\alpha - \boldsymbol{\epsilon} = (\mathbf{a}^\alpha \otimes \mathbf{n})_S$.

Unlike for the strains, the recovery of the phase-specific stresses is most easily done along the same lines as for the total stress. Starting again in the \mathcal{B} -system, the normal stress components are given by $\boldsymbol{\sigma}_n^\alpha = \boldsymbol{\sigma}_n$, whereas the tangential ones can either be obtained directly from the phase-specific strains calculated above, or, more in line with the previous calculations, through

$$\begin{aligned} \boldsymbol{\sigma}_t^\alpha &= \mathbf{C}_{tn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) = \mathbf{C}_{tn}^\alpha (\boldsymbol{\epsilon}_n - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) + \mathbf{C}_{tn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \boldsymbol{\epsilon}_n) \\ &= \boldsymbol{\Sigma}_t^\alpha + \mathbf{C}_{tn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \boldsymbol{\epsilon}_n) = \boldsymbol{\Sigma}_t^\alpha + \mathbf{C}_{tn}^\alpha \mathbf{S}_{nn}^\alpha (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^\alpha). \end{aligned}$$

The equality $\boldsymbol{\sigma}_n^\alpha = \boldsymbol{\sigma}_n$ can again, similar to its phase-averaged version in Equation (7.116), artificially be expanded to $\boldsymbol{\sigma}_n^\alpha = \boldsymbol{\Sigma}_n^\alpha + \mathbf{C}_{nn}^\alpha \mathbf{S}_{nn}^\alpha (\boldsymbol{\sigma}_n - \boldsymbol{\Sigma}_n^\alpha)$. Based on this common representation of the full stress tensor, the same arguments as for the retransformation of $\boldsymbol{\sigma}_{\mathcal{B}}$ can be applied, thus allowing to express the entire phase-specific stress tensor in the original system through

$$\boldsymbol{\sigma}^\alpha = \boldsymbol{\Sigma}^\alpha + (\mathbf{C}^\alpha \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \left((\boldsymbol{\sigma} - \boldsymbol{\Sigma}^\alpha) \cdot \mathbf{n} \right). \quad (7.130)$$

Remark 141. In relation with Remark 138, it should be noted that the quite straightforward retransformation of the normal strain-jumps in contrast to the more difficult “retransformation” of the normal stress components is due to the interplay of the additional factor 2 on the shear-components in the Voigt-representation of the strain-type tensors in combination with Equation

(7.103). In fact, as the strain-jump and the “extension by zero” $\begin{pmatrix} \sigma_{nn} & \sigma_{nt} \\ \sigma_{nt} & 0 \end{pmatrix}$ of the normal

stress have the same structure, it is clear that this extension could also be written in terms of a symmetrized dyad $(\mathbf{b} \otimes \mathbf{n})_S$ with some vector \mathbf{b} . By Equation (7.103), this vector \mathbf{b} is given by $\mathbf{b} = (\sigma_{nn}, 2\sigma_{nt})^T$ (this multiplication of the shear-component by 2 being the role of the rightmost-matrix in Equation (7.123)). The middle factor \mathbf{Q}^T then simply serves to retransform this vector back to the normal coordinate system. Finally, the left-most matrix - corresponding to the \mathbf{B} -matrix with the “shear-entries halved” - is the Voigt-representation of a symmetrized dyad with \mathbf{n} for stress-type tensors in Voigt notation, i.e. one where the shear-components are not stored in doubled form.

That the retransformation of the normal strain jump is so easily expressed in term of \mathbf{B}^T and \mathbf{Q}^T is therefore somewhat fortuitous. On the one hand, given that the shear-components of the strain are stored in doubled form, the 2-factor in the calculation of the jump-vector \mathbf{a} according to Equation (7.103) is already included. On the other hand, storing the symmetrized dyad $(\mathbf{a} \otimes \mathbf{n})_S$ with doubled shear-components cancels the $\frac{1}{2}$ -prefactor in the left-most matrix in Equation (7.123), leading to \mathbf{B}^T . \diamond

Remark 142. From an implementation point of view, at least in the current setting, the (obviously equivalent) alternative expression

$$\boldsymbol{\sigma} = \sum_{\alpha=1}^N \left(\boldsymbol{\Sigma}^\alpha - (\boldsymbol{C}^\alpha \cdot \boldsymbol{n}) \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) \right) h^\alpha(\phi) + \left(\sum_{\alpha=1}^N (\boldsymbol{C}^\alpha \cdot \boldsymbol{n}) \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha(\phi) \right) \cdot (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \quad (7.131)$$

for the stress is in principle more convenient than the one in Equation (7.125). The advantage of this form in combination with the expression (repeated from Equation (7.122) for convenience)

$$\boldsymbol{\sigma} \cdot \boldsymbol{n} = \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha(\phi) \right)^{-1} \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) h^\alpha(\phi)$$

for the normal stress is that the first term contains only initially known quantities, whereas the second part consist out of a **single** correction operation based on an averaged \boldsymbol{C}^α - \boldsymbol{n} -combination in the end once $\boldsymbol{\sigma} \cdot \boldsymbol{n}$ has been calculated. In addition, all phase-specific quantities required for both calculating $\boldsymbol{\sigma} \cdot \boldsymbol{n}$ as well as the correction are essentially the same as the ones already required for the evaluation of the first term. It is therefore not necessary to evaluate these two times - once for the determination of $\boldsymbol{\sigma} \cdot \boldsymbol{n}$ and then for the evaluation of $\boldsymbol{\sigma}$ itself once the normal stress is known - as one can simply accumulate the required quantities while running over the phases in the evaluation of the $\boldsymbol{\Sigma}^\alpha - (\boldsymbol{C}^\alpha \cdot \boldsymbol{n}) \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n})$.

For more complex models where one requires access to the phase-specific stresses and/or strains for the stress-calculation itself (such as for visco-elastic or plastic problems), an obvious alternative to both equations (7.125) and (7.131) is to simply calculate the total stress as the average of the phase-specific ones (presumably already calculated based upon e.g. Equation (7.130)), $\boldsymbol{\sigma} = \sum_{\alpha} \boldsymbol{\sigma}^\alpha h^\alpha(\phi)$.

◇

Remark 143. In either case, it should be kept in mind that regardless of the particular formulation and implementation used, the jump-condition based models are fairly expensive to evaluate within the interface region (this being even more problematic for potential multiphase extensions). Whenever - such as in many implicit or linearized schemes - a larger amount of evaluations of stress-increments have to be performed based on fixed phasefield values (an potentially other parameters), it may be significantly cheaper from a computational point of view to precalculate the “effective” or algorithmically consistent stiffness, linking the increments in the total stress to those of the total strain. For the linearly elastic case above, it is straightforward to verify that, given fixed eigenstrains $\tilde{\boldsymbol{\epsilon}}^\alpha$, an increment in the average strain $\boldsymbol{\epsilon}$ leads to the increment

$$\delta(\boldsymbol{\sigma} \cdot \boldsymbol{n}) = \left(\sum_{\alpha=1}^N (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha(\phi) \right)^{-1} \cdot \left(\sum_{\alpha=1}^N (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha) h^\alpha(\phi) \right) : \delta \boldsymbol{\epsilon}.$$

Combined with the increments $\delta \boldsymbol{\Sigma}^\alpha = \boldsymbol{C}^\alpha : \delta \boldsymbol{\epsilon}$ for the Voigt-Taylor stress prediction, it then follows from Equation (7.125) that the increment in the effective stress $\delta \boldsymbol{\sigma}$ is related to the one in the total strain by $\delta \boldsymbol{\sigma} = \boldsymbol{C}_{eff} : \delta \boldsymbol{\epsilon}$ with \boldsymbol{C}_{eff} given by

$$\boldsymbol{C}_{eff} = \sum_{\alpha=1}^2 \left[\boldsymbol{C}^\alpha + (\boldsymbol{C}^\alpha \cdot \boldsymbol{n}) \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot \left(\left(\sum_{\beta=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\beta \cdot \boldsymbol{n})^{-1} h^\beta \right)^{-1} \cdot \left(\sum_{\beta=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\beta \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{n} \cdot \boldsymbol{C}^\beta) h^\beta \right) - \boldsymbol{n} \cdot \boldsymbol{C}^\alpha \right) \right] h^\alpha. \quad (7.132)$$

◇

As a final step, it remains to derive an expression for the driving force in the original coordinate system. The simplest way of doing so is simply reusing the expression

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\beta=1}^2 \left(f_{el}^{\beta} - \boldsymbol{\sigma}_n \cdot \boldsymbol{\epsilon}_n^{\beta} \right) \frac{\partial h^{\beta}}{\partial \phi} = \sum_{\beta=1}^2 \left(\frac{1}{2} (\boldsymbol{\epsilon}_B^{\beta} - \tilde{\boldsymbol{\epsilon}}_B^{\beta}) \cdot \mathbf{C}_B^{\beta} (\boldsymbol{\epsilon}_B^{\beta} - \tilde{\boldsymbol{\epsilon}}_B^{\beta}) - \boldsymbol{\sigma}_n \cdot \boldsymbol{\epsilon}_n^{\beta} \right) \frac{\partial h^{\beta}}{\partial \phi}$$

derived in [64]. The first part f_{el}^{β} can directly be reexpressed (the transformation matrices for the strains and the stiffness tensor simply cancel) using either the expressions for $(\boldsymbol{\epsilon}^v)^{\alpha}$ in Equation (7.127) together with the original stiffness-tensor, or by combining the expression for $\boldsymbol{\epsilon}^{\alpha}$ with the one for the phase-specific stresses $(\boldsymbol{\sigma}^v)^{\alpha}$ from Equation (7.130).

Due to its particular form, the additional contribution by the normal stress can also easily be rewritten in terms of quantities in the original system. In fact, firstly using

$$\sum_{\beta=1}^2 \boldsymbol{\sigma}_n \cdot \boldsymbol{\epsilon}_n^{\beta} \frac{\partial h^{\beta}}{\partial \phi} = \sum_{\beta=1}^2 \boldsymbol{\sigma}_n \cdot (\boldsymbol{\epsilon}_n^{\beta} - \boldsymbol{\epsilon}_n) \frac{\partial h^{\beta}}{\partial \phi} + \underbrace{\boldsymbol{\sigma}_n \cdot \boldsymbol{\epsilon}_n}_{=0} \sum_{\beta=1}^2 \frac{\partial h^{\beta}}{\partial \phi}$$

and combining this with $\boldsymbol{\sigma}_n = \mathbf{Q} \mathbf{B} \boldsymbol{\sigma}^v$ and the expression (7.126) for the phase-specific strains and the definition $\boldsymbol{\Sigma}_n^{\alpha} = \mathbf{Q} \mathbf{B} (\boldsymbol{\Sigma}^v)^{\alpha}$ leads to

$$\boldsymbol{\sigma}_n \cdot \boldsymbol{\epsilon}_n^{\beta} = \mathbf{Q} \mathbf{B} \boldsymbol{\sigma}^v \cdot \mathbf{S}_{nn}^{\alpha} \mathbf{Q} \mathbf{B} (\boldsymbol{\sigma}^v - \boldsymbol{\Sigma}_n^{\alpha}) = \mathbf{B} \boldsymbol{\sigma}^v \cdot \mathbf{Q}^T \mathbf{S}_{nn}^{\alpha} \mathbf{Q} \mathbf{B} (\boldsymbol{\sigma}^v - \boldsymbol{\Sigma}^{\alpha}).$$

The \mathbf{Q} -factors cancel with the ones in the expression $\mathbf{S}_{nn}^{\beta} = \mathbf{Q} (\mathbf{B} (\mathbf{C}^v)^{\beta} \mathbf{B}^T)^{-1} \mathbf{Q}^T$ from Equation (7.119) though, thus leaving

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\beta=1}^2 \left(f_{el}^{\beta} - (\mathbf{B} \boldsymbol{\sigma}^v) \cdot (\mathbf{B} (\mathbf{C}^v)^{\beta} \mathbf{B}^T \mathbf{B} (\boldsymbol{\sigma}^v - \boldsymbol{\Sigma}^{\alpha})) \right) \frac{\partial h^{\beta}}{\partial \phi}, \quad (7.133)$$

or, in non-Voigt notation,

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\beta=1}^2 \left(f_{el}^{\beta} - (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \mathbf{C}^{\beta} \cdot \mathbf{n}) \cdot ((\boldsymbol{\sigma} - \boldsymbol{\Sigma}^{\alpha}) \cdot \mathbf{n}) \right) \frac{\partial h^{\beta}}{\partial \phi}. \quad (7.134)$$

Remark 144. Note that, with the definition (7.128) of \mathbf{a}^{α} , these expressions can also be rewritten in the more pleasant form

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\beta=1}^2 \left(f_{el}^{\beta} - (\mathbf{B} \boldsymbol{\sigma}^v) \cdot \mathbf{a}^{\alpha} \right) \frac{\partial h^{\beta}}{\partial \phi} \quad \text{resp.} \quad \frac{\partial f_{el}}{\partial \phi} = \sum_{\beta=1}^2 \left(f_{el}^{\beta} - (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{a}^{\alpha} \right) \frac{\partial h^{\beta}}{\partial \phi}. \quad (7.135)$$

With $\mathbf{B} \boldsymbol{\sigma}^v \cdot \mathbf{a}^{\alpha} = \boldsymbol{\sigma}^v \cdot (\mathbf{B}^T \mathbf{a}^{\alpha})$ (resp. $(\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{a}^{\alpha} = \boldsymbol{\sigma} : (\mathbf{n} \otimes \mathbf{a})_S$) using the symmetry of $\boldsymbol{\sigma}$ ⁷³, it follows that yet another form of the driving force is given by

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\alpha=1}^2 \left(f_{el}^{\alpha} - \boldsymbol{\sigma} : (\mathbf{a}^{\alpha} \otimes \mathbf{n})_S \right) \frac{\partial h^{\alpha}}{\partial \phi}, \quad (7.136)$$

which, as $(\mathbf{a}^{\alpha} \otimes \mathbf{n})_S = \boldsymbol{\epsilon}^{\alpha} - \boldsymbol{\epsilon}$ and $\sum_{\alpha=1}^2 \boldsymbol{\sigma} : \boldsymbol{\epsilon} \frac{\partial h^{\alpha}}{\partial \phi} = \boldsymbol{\sigma} : \boldsymbol{\epsilon} \left(\sum_{\alpha=1}^2 \frac{\partial h^{\alpha}}{\partial \phi} \right) = 0$, can finally also be replaced by

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\alpha=1}^2 \left(f_{el}^{\alpha} - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^{\alpha} \right) \frac{\partial h^{\alpha}}{\partial \phi}. \quad (7.137)$$

◇

⁷³As in Equation (7.129), $\mathbf{B}^T \mathbf{a}^{\alpha}$ is precisely the Voigt-representation of the (strain-type) tensor $(\mathbf{a} \otimes \mathbf{n})_S$.

Remark 145. It is also possible to directly derive the (same) driving force directly based upon expression in the original coordinate system only. In fact, from $\boldsymbol{\sigma}^\alpha = \frac{\partial f^\alpha}{\partial \boldsymbol{\epsilon}^\alpha}$, one has

$$\frac{\partial f_{el}}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{\alpha=1}^2 f^\alpha h^\alpha = \sum_{\alpha=1}^2 f^\alpha \frac{\partial h^\alpha}{\partial \phi} + \boldsymbol{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) : \frac{\partial \boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \phi} h^\alpha = \sum_{\alpha=1}^2 f^\alpha \frac{\partial h^\alpha}{\partial \phi} + \boldsymbol{\sigma}^\alpha : \frac{\partial \boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \phi} h^\alpha. \quad (7.138)$$

As neither $\tilde{\boldsymbol{\epsilon}}^\alpha$ nor $\boldsymbol{\epsilon}$ depend on ϕ and $\boldsymbol{\epsilon}^\alpha - \boldsymbol{\epsilon} = (\boldsymbol{a}^\alpha \otimes \boldsymbol{n})_S$ and $\boldsymbol{\sigma}^\alpha$ is symmetric, the second contribution can be rewritten as

$$\sum_{\alpha=1}^2 \boldsymbol{\sigma}^\alpha : \frac{\partial \boldsymbol{\epsilon}^\alpha - \boldsymbol{\epsilon}}{\partial \phi} h^\alpha = \sum_{\alpha=1}^2 \boldsymbol{\sigma}^\alpha : \left(\frac{\partial \boldsymbol{a}^\alpha}{\partial \phi} \otimes \boldsymbol{n} \right)_S h^\alpha = \sum_{\alpha=1}^2 (\boldsymbol{\sigma}^\alpha \cdot \boldsymbol{n}) \cdot \frac{\partial \boldsymbol{a}^\alpha}{\partial \phi} h^\alpha.$$

By construction, $\boldsymbol{\sigma}^\alpha \cdot \boldsymbol{n} = \boldsymbol{\sigma} \cdot \boldsymbol{n}$, such that it only remains to evaluate $\frac{\partial \boldsymbol{a}^\alpha}{\partial \phi}$. By Equation (7.128), it further follows that

$$\frac{\partial \boldsymbol{a}^\alpha}{\partial \phi} = (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot \frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi}$$

and thus that

$$\sum_{\alpha=1}^2 \boldsymbol{\sigma}^\alpha : \frac{\partial \boldsymbol{\epsilon}^\alpha - \boldsymbol{\epsilon}}{\partial \phi} h^\alpha = (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right) \cdot \frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi}. \quad (7.139)$$

From Equation (7.122), one has

$$\begin{aligned} \frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi} &= \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) \frac{\partial h^\alpha}{\partial \phi} \\ &\quad + \left(\frac{\partial}{\partial \phi} \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \right) \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) h^\alpha, \end{aligned}$$

where the derivative in the second term can be evaluated as

$$-\left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \frac{\partial h^\alpha}{\partial \phi} \right) \cdot \underbrace{\left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) h^\alpha}_{=\boldsymbol{\sigma} \cdot \boldsymbol{n}}.$$

Reinserting this expression into the derivative of $\boldsymbol{\sigma} \cdot \boldsymbol{n}$ then leads to

$$\begin{aligned} \frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi} &= \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) \frac{\partial h^\alpha}{\partial \phi} \\ &\quad - \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \frac{\partial h^\alpha}{\partial \phi} \right) \cdot (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \\ &= \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^\alpha \cdot \boldsymbol{n}) - (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \right) \frac{\partial h^\alpha}{\partial \phi}. \end{aligned}$$

Recognizing that the summands in the second sum correspond to \boldsymbol{a}^α , this further simplifies to

$$\frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi} = - \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right)^{-1} \cdot \sum_{\alpha=1}^2 \boldsymbol{a}^\alpha \frac{\partial h^\alpha}{\partial \phi}, \quad (7.140)$$

and finally, combined with Equation (7.138) and (7.139), to

$$\frac{\partial f_{el}}{\partial \phi} = \sum_{\alpha=1}^2 f^\alpha \frac{\partial h^\alpha}{\partial \phi} + (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \cdot \left(\sum_{\alpha=1}^2 (\boldsymbol{n} \cdot \boldsymbol{C}^\alpha \cdot \boldsymbol{n})^{-1} h^\alpha \right) \cdot \frac{\partial (\boldsymbol{\sigma} \cdot \boldsymbol{n})}{\partial \phi} = \sum_{\alpha=1}^2 (f^\alpha - (\boldsymbol{\sigma} \cdot \boldsymbol{n}) \cdot \boldsymbol{a}^\alpha) \frac{\partial h^\alpha}{\partial \phi}.$$

This is precisely the same expression as the second one in Equation (7.135). \diamond

In contrast to the (correct) derivative with respect to ϕ , there is no expression for the derivative with respect to $\nabla\phi$ in [64], as the paper (falsely) claims that this derivative vanishes. In addition, given that the transformation matrices \mathbf{M}_σ and \mathbf{M}_ϵ depend in an intricate fashion on the normal and (in particular in the three-dimensional case) how the tangent vector(s) are constructed based upon \mathbf{n} , performing this derivative based upon the formulation in [64] is likely to be a very arduous task. In contrast, it is straightforward based upon the reformulation above and runs essentially along the same lines - but actually in a simpler form - as the derivation of the expression for $\frac{\partial f_{el}}{\partial\phi}$ in the previous remark.

Again starting from $f_{el} = \sum_{\alpha=1}^2 \frac{1}{2}(\epsilon^\alpha - \tilde{\epsilon}^\alpha) : \mathbf{C}^\alpha : (\epsilon^\alpha - \tilde{\epsilon}^\alpha) h^\alpha(\phi)$, a differentiation w.r.t. \mathbf{n} leads to

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = \sum_{\alpha=1}^2 \boldsymbol{\sigma}^\alpha : \frac{\partial(\epsilon^\alpha - \tilde{\epsilon}^\alpha)}{\partial \mathbf{n}} h^\alpha = \sum_{\alpha=1}^2 \boldsymbol{\sigma}^\alpha : \frac{\partial}{\partial \mathbf{n}} (\mathbf{a}^\alpha \otimes \mathbf{n})_S h^\alpha$$

since, similar to above, $\frac{\partial(\epsilon^\alpha - \tilde{\epsilon}^\alpha)}{\partial \mathbf{n}} = \frac{\partial(\epsilon^\alpha - \epsilon)}{\partial \mathbf{n}}$ as both $\tilde{\epsilon}^\alpha$ and ϵ are independent of \mathbf{n} . By the symmetry of $\boldsymbol{\sigma}$, this can be simplified to

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = \sum_{\alpha=1}^2 \left((\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}^\alpha}{\partial \mathbf{n}} + (\boldsymbol{\sigma}^\alpha \cdot \mathbf{a}^\alpha) \cdot \frac{\partial \mathbf{n}}{\partial \mathbf{n}} \right) h^\alpha = \sum_{\alpha=1}^2 \left((\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}^\alpha}{\partial \mathbf{n}} + (\boldsymbol{\sigma}^\alpha \cdot \mathbf{a}^\alpha) \right) h^\alpha,$$

where use was made of the equality of the phase-specific normal stresses with the total one. As $\boldsymbol{\sigma} \cdot \mathbf{n}$ does not depend on α , the first contribution can be summarized to $(\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \sum_{\alpha=1}^2 \frac{\partial(\sum_{\alpha=1}^2 \mathbf{a}^\alpha h^\alpha)}{\partial \mathbf{n}}$, and since $\sum_{\alpha=1}^2 \mathbf{a}^\alpha h^\alpha = \mathbf{0}$ ⁷⁴, actually drops out, leaving the very simple expression

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = \sum_{\alpha=1}^2 (\boldsymbol{\sigma}^\alpha \cdot \mathbf{a}^\alpha) h^\alpha \quad (7.141)$$

for $\frac{\partial f_{el}}{\partial \mathbf{n}}$.

Remark 146. An analogous calculation to the one above could also be performed in the three-dimensional case. The only major difference to the two-dimensional one is the larger dimensionality of the vectors and subvectors as the Voigt-representation then consists of six-dimensional vectors and both the normal and tangential components are made up of three entries. For example, the \mathbf{M}_σ -matrix for performing the transformation to the \mathcal{B} -system becomes [64]

$$\mathbf{M}_\sigma = \begin{pmatrix} n_x^2 & n_y^2 & n_z^2 & 2n_y n_z & 2n_x n_z & 2n_x n_y \\ t_x^2 & t_y^2 & t_z^2 & 2t_y t_z & 2t_x t_z & 2t_x t_y \\ s_x^2 & s_y^2 & s_z^2 & 2s_y s_z & 2s_x s_z & 2s_x s_y \\ t_x s_x & t_y s_y & t_z s_z & t_y s_z + t_z s_y & t_x s_z + t_z s_x & t_x s_y + t_y s_x \\ n_x s_x & n_y s_y & n_z s_z & n_y s_z + n_z s_y & n_x s_z + n_z s_x & n_x s_y + n_y s_x \\ n_x t_x & n_y t_y & n_z t_z & n_y t_z + n_z t_y & n_x t_z + n_z t_x & n_x t_y + n_y t_x \end{pmatrix}.$$

With this transformation matrix, reading out the normal components in the \mathcal{B} -system (in the order $(\boldsymbol{\sigma}_{nn}, \boldsymbol{\sigma}_{nt}, \boldsymbol{\sigma}_{ns})^T$) is represented by the combined matrix

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \mathbf{M}_\sigma = \begin{pmatrix} n_x^2 & n_y^2 & n_z^2 & 2n_y n_z & 2n_x n_z & 2n_x n_y \\ n_x t_x & n_y t_y & n_z t_z & n_y t_z + n_z t_y & n_x t_z + n_z t_x & n_x t_y + n_y t_x \\ n_x s_x & n_y s_y & n_z s_z & n_y s_z + n_z s_y & n_x s_z + n_z s_x & n_x s_y + n_y s_x \end{pmatrix} \\ = \underbrace{\begin{pmatrix} n_x & n_y & n_z \\ t_x & t_y & t_z \\ s_x & s_y & s_z \end{pmatrix}}_{=: \mathbf{Q}} \underbrace{\begin{pmatrix} n_x & 0 & 0 & 0 & n_z & n_y \\ 0 & n_y & 0 & n_z & 0 & n_x \\ 0 & 0 & n_z & n_y & n_x & 0 \end{pmatrix}}_{=: \mathbf{B}},$$

⁷⁴This can be verified “manually” based on Equation (7.128) combined with the definition of $\boldsymbol{\sigma} \cdot \mathbf{n}$, but is also easily seen from the imposed averaging relation $\mathbf{0} = \sum_{\alpha=1}^2 \epsilon^\alpha h^\alpha - \epsilon = \sum_{\alpha=1}^2 (\epsilon^\alpha - \epsilon) h^\alpha = \sum_{\alpha=1}^2 \mathbf{a}^\alpha h^\alpha$.

where \mathbf{Q} is again a unitary matrix corresponding to the transformation of “standard” vectors, and \mathbf{B} encodes the multiplication of a Voigt-type quantity with \mathbf{n} . Adjusting the remaining formulae above making explicit use of the two-dimensional nature in an analogous manner, one can then verify that the rest of the calculation can be performed along exactly the same lines, just requiring more “bulky” matrices. \diamond

Comparison with the Model by Durga et al. [24]

As already indicated in Section 7.2.2, Durga et al. propose a slightly different model in [24] (generalizing their previous work [23]). The formulation is in fact a priori very similar to the one in [64] in the sense that both approaches rely upon a transformation into a coordinate system where the normal vector \mathbf{n} corresponds to the first unit vector \mathbf{e}_1 and that, within this transformed system and using the notation from [64], one assumes $\boldsymbol{\sigma}_n^\alpha = \boldsymbol{\sigma}_n$, $\boldsymbol{\epsilon}_t^\alpha = \boldsymbol{\epsilon}_t$ and $\boldsymbol{\sigma}_t^\alpha = \mathbf{C}_{tn}^\alpha(\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t)$ ⁷⁵. In contrast to the model in [64], the (common) normal stress is a priori not (but see the discussion below) calculated from $\boldsymbol{\sigma}_n^\alpha = \mathbf{C}_{nn}^\alpha(\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{nt}^\alpha(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t)$, but instead simply assumed to be equal to the corresponding components of the one calculated using the Reuss-Sachs model. In non-Voigt notation, one therefore has

$$\boldsymbol{\sigma} \cdot \mathbf{n} := \boldsymbol{\sigma}_{RS} \cdot \mathbf{n} = \left(\left(\sum_{\alpha=1}^2 \mathbf{S}^\alpha h^\alpha(\phi) \right)^{-1} : \left(\boldsymbol{\epsilon} - \sum_{\alpha=1}^2 \tilde{\boldsymbol{\epsilon}}^\alpha h^\alpha(\phi) \right) \right) \cdot \mathbf{n}.$$

In order to complete the model, it remains to specify a rule for determining the normal strains $\boldsymbol{\epsilon}_n^\alpha$. This is achieved here - in line with the argument underlying the Reuss-Sachs model - by defining $\boldsymbol{\epsilon}_n^\alpha$ to be given by the normal components of $\mathbf{S}^\alpha : \boldsymbol{\sigma}^\alpha$, i.e. $\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha = (\mathbf{S}^\alpha)_{nn} \boldsymbol{\sigma}_n + (\mathbf{S}^\alpha)_{nt} \boldsymbol{\sigma}_t^\alpha$, where $\boldsymbol{\sigma}^\alpha$ is the (full) stress reconstructed based upon the normal components from the Reuss-Sachs model and $\boldsymbol{\sigma}_t^\alpha$ based upon the phase-specific material law⁷⁶. Combined with the expression for $\boldsymbol{\sigma}_t^\alpha$ in terms of $\boldsymbol{\epsilon}_n^\alpha$ leads to

$$\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha = (\mathbf{S}^\alpha)_{nn} \boldsymbol{\sigma}_n + (\mathbf{S}^\alpha)_{nt} \left(\mathbf{C}_{tn}^\alpha(\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha(\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) \right),$$

resp.

$$\left(\mathbf{I} - (\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tn}^\alpha \right) (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) = (\mathbf{S}^\alpha)_{nn} \boldsymbol{\sigma}_n + (\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) \quad (7.142)$$

and thus to an equation which can be solved for $\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha$ given the prescribed normal stress and tangential strains.

This equality can actually be modified into a more instructive form as \mathbf{S}^α and \mathbf{C}^α are inverses (in Voigt notation) of each other, i.e. one has

$$\begin{pmatrix} (\mathbf{S}^\alpha)_{nn} & (\mathbf{S}^\alpha)_{nt} \\ (\mathbf{S}^\alpha)_{tn} & (\mathbf{S}^\alpha)_{tt} \end{pmatrix} \begin{pmatrix} \mathbf{C}_{nn}^\alpha & \mathbf{C}_{nt}^\alpha \\ \mathbf{C}_{tn}^\alpha & \mathbf{C}_{tt}^\alpha \end{pmatrix} = \begin{pmatrix} (\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nn}^\alpha + (\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tn}^\alpha & (\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nt}^\alpha + (\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tt}^\alpha \\ (\mathbf{S}^\alpha)_{tn} \mathbf{C}_{nn}^\alpha + (\mathbf{S}^\alpha)_{tt} \mathbf{C}_{tn}^\alpha & (\mathbf{S}^\alpha)_{tn} \mathbf{C}_{nt}^\alpha + (\mathbf{S}^\alpha)_{tt} \mathbf{C}_{tt}^\alpha \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}.$$

From the upper-left block of this equality, one has $\mathbf{I} - (\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tn}^\alpha = (\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nn}^\alpha$, while the upper-right block shows that $(\mathbf{S}^\alpha)_{nt} \mathbf{C}_{tt}^\alpha = -(\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nt}^\alpha$. Inserting these expressions into Equation (7.142), the phase-specific normal strains can alternatively be characterized by

$$(\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) = (\mathbf{S}^\alpha)_{nn} \boldsymbol{\sigma}_n - (\mathbf{S}^\alpha)_{nn} \mathbf{C}_{nt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t),$$

and thus, $(\mathbf{S}^\alpha)_{nn}$ being invertible as a diagonal subblock of a positive definite matrix, in fact

$$\boldsymbol{\sigma}_n = \mathbf{C}_{nn}^\alpha (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{nt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t). \quad (7.143)$$

⁷⁵In the Voigt-index based notation of the appendix A in [24], the normal and tangential stress components are given by $\boldsymbol{\sigma}_n = (\sigma_1, \sigma_5, \sigma_6)^T$ and $\boldsymbol{\sigma}_t = (\sigma_2, \sigma_3, \sigma_4)^T$ (and similarly for the strains).

⁷⁶Note the difference in notation here with $(\mathbf{S}^\alpha)_{nn}$ designating the actual $\mathbf{n} - \mathbf{n}$ -subblock of \mathbf{S}^α , in contrast with \mathbf{S}_{nn}^α defined above.

From this, it is obvious that the only difference between the model in [24] and the one discussed above⁷⁷ is that in the former, the normal stresses are prescribed directly based on the Reuss-Sachs/Steinbach-Apel scheme, whereas in the latter case, $\boldsymbol{\sigma}_n$ is obtained by enforcing the averaging condition

$$\boldsymbol{\epsilon}_n \stackrel{!}{=} \sum_{\alpha=1}^2 \boldsymbol{\epsilon}_n^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^2 (\mathbf{C}_{nn}^\alpha)^{-1} (\boldsymbol{\sigma}_n - \mathbf{C}_{nt}^\alpha (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha)) h^\alpha(\phi)$$

based upon Equation (7.143).

In contrast to the model by [64], the model in [24] therefore does not generally satisfy the natural relation $\sum_{\alpha=1}^2 \boldsymbol{\epsilon}_n^\alpha h^\alpha(\phi) = \boldsymbol{\epsilon}_n$ unless the normal stresses in e.g. Equation (7.108) or Equation (7.122) happen to coincide with⁷⁸ $\boldsymbol{\sigma}_{RS} \cdot \mathbf{n}$.

Remark 147. As the Reuss-Sachs/Steinbach-Apel model is based on the assumption of (fully) equal stresses between all phases, it is, as used in [23] and [24], an intuitively much more reasonable choice for a common normal stress than the one from e.g. the Voigt-Taylor model, which is based on a completely different assumption. In addition, as shown in these papers, combining this assumption with the subsequent imposition of common tangential strains does in fact lead to a marked improvement in comparison with both the Reuss-Sachs- and Voigt-Taylor-model themselves.

Nevertheless, the stress-prediction by the Reuss-Sachs model is (implicitly) based on phase-specific strains where the tangential components need not coincide. In combination with Equation (7.143), the subsequent choice of the phase-specific normal strains based on $\boldsymbol{\sigma}_{RS} \cdot \mathbf{n}$ can then be interpreted as adjusting the normal strains such that, now under the different (and more physical assumption) $\boldsymbol{\epsilon}_t^\alpha = \boldsymbol{\epsilon}_t$ and with the phase-specific material law, one obtains this same “predicted” normal stress for both phases. Even though this model presents a major step forward in terms of modeling the stress-strain relationship within a diffuse interface, it does therefore not achieve the same degree of internal consistency as the model(s) by [51] and [64], where the calculation of $\boldsymbol{\sigma} \cdot \mathbf{n}$ itself is already based on the same assumptions as the evaluation of the final stresses⁷⁹. \diamond

Equivalence of the Two-Phase Model to the One by Tschukin

Yet another formulation for the jump-condition based mechanical model is presented by Tschukin in [74]. As the primary motivations of this alternative formulation to the model in [64] (a collaboration between, among others, Schneider with Tschukin), Tschukin also mentions the computational overhead involved in the transformations to the \mathcal{B} -system (including the necessity of generating two tangent vectors instead of relying solely on \mathbf{n}) as well as correcting the wrong claim of the vanishing contribution $\frac{\partial f_{\epsilon_t}}{\partial \nabla \phi}$ in the variational derivative.

Even though the description in [74] is in many ways formally similar to the one in [64] in the sense that it continues to decompose both the strain- and stress-tensor (and similarly the stiffness tensor) into normal and tangential “components”, the crucial difference is that these are not defined in terms of a collection of specific entries within the transformed coordinate system \mathcal{B} anymore. Instead, Tschukin defines these quantities in terms of two (forth-order) projection operators

$$\mathcal{N} = \mathbf{I} \square \boldsymbol{\Phi} + \boldsymbol{\Phi} \square \mathbf{I} - \boldsymbol{\Phi} \square \boldsymbol{\Phi} \quad \text{and} \quad \mathcal{T} = (\mathbf{I} - \boldsymbol{\Phi}) \square (\mathbf{I} - \boldsymbol{\Phi}) \quad (7.144)$$

⁷⁷As one explicitly assumes $\boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon}_t$, the averaging condition $\boldsymbol{\epsilon}_t = \sum_{\alpha=1}^2 \boldsymbol{\epsilon}_t^\alpha h^\alpha = \boldsymbol{\epsilon}_t$ is trivially satisfied, and from Equation (7.143) combined with the analogous (explicit) assumption for the tangential stress, the model also satisfies the phase-specific material law $\boldsymbol{\sigma}^\alpha = \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$, from which the total stress is obtained by a weighted interpolation (trivially reducing to $\boldsymbol{\sigma}_n$ for the normal components).

⁷⁸One (quite particular) case where this actually happens is when both materials are taken as isotropic, with the Lamé-parameters of one phase being chosen as (equal) multiples of those of the other phase, i.e. $\mathbf{C}^2 = c\mathbf{C}^1$ with $\mathbf{C}^1 = \lambda^1 \mathbf{I} \otimes \mathbf{I} + 2\mu^1 \mathbf{I}^{(s)}$.

⁷⁹Note also, that as in [64], the dependence on \mathbf{n} is disregarded in the phasefield evolution equation in [23] and [24].

onto the normal and tangential components, where $\Phi = \mathbf{n} \otimes \mathbf{n}$ and the “box-product” is defined by $(\mathbf{A} \square \mathbf{B}) : \mathbf{C} = \mathbf{A} \cdot \mathbf{C} \cdot \mathbf{B}$. In terms of these, the normal and tangential components of the stress- and strain-tensor are defined by

$$\boldsymbol{\sigma}_n := \mathcal{N} : \boldsymbol{\sigma}, \quad \boldsymbol{\sigma}_t := \mathcal{T} : \boldsymbol{\sigma}, \quad \text{and} \quad \boldsymbol{\epsilon}_n := \mathcal{N} : \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon}_t := \mathcal{T} : \boldsymbol{\epsilon}. \quad (7.145)$$

Using $\boldsymbol{\sigma}^\alpha = \mathcal{N} : \boldsymbol{\sigma}^\alpha + \mathcal{T} : \boldsymbol{\sigma}^\alpha$ and $\boldsymbol{\epsilon}^\alpha = \mathcal{N} : \boldsymbol{\epsilon}^\alpha + \mathcal{T} : \boldsymbol{\epsilon}^\alpha$ together with the idempotence and orthogonality of these projectors, this allows rewriting the normal “part” of the phase-specific material law as

$$\begin{aligned} \boldsymbol{\sigma}_n^\alpha &= \mathcal{N} : \boldsymbol{\sigma}^\alpha = \mathcal{N} : \mathbf{C}^\alpha : (\mathcal{N} : \boldsymbol{\epsilon}^\alpha + \mathcal{T} : \boldsymbol{\epsilon}^\alpha) = \mathcal{N} : \mathbf{C}^\alpha : \mathcal{N} : \boldsymbol{\epsilon}^\alpha + \mathcal{N} : \mathbf{C}^\alpha : \mathcal{T} : \boldsymbol{\epsilon}^\alpha \\ &= (\mathcal{N} : \mathbf{C}^\alpha : \mathcal{N}) : (\mathcal{N} : \boldsymbol{\epsilon}^\alpha) + \mathcal{N} : \mathbf{C}^\alpha : \mathcal{T} : (\mathcal{T} : \boldsymbol{\epsilon}^\alpha) = (\mathcal{N} : \mathbf{C}^\alpha : \mathcal{N}) : \boldsymbol{\epsilon}_n^\alpha + \mathcal{N} : \mathbf{C}^\alpha : \mathcal{T} : \boldsymbol{\epsilon}_t^\alpha. \end{aligned}$$

Proceeding similarly with the tangential part and adding, by the same argument, the eigenstrains, the material behavior can be fully characterized by $\boldsymbol{\sigma}^\alpha = \boldsymbol{\sigma}_n^\alpha + \boldsymbol{\sigma}_t^\alpha$ through

$$\begin{aligned} \boldsymbol{\sigma}_n^\alpha &= \mathbf{C}_{nn}^\alpha : (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon}_t^\alpha - \tilde{\boldsymbol{\epsilon}}_t^\alpha) \quad \text{and} \\ \boldsymbol{\sigma}_t^\alpha &= \mathbf{C}_{tn}^\alpha : (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha : (\boldsymbol{\epsilon}_t^\alpha - \tilde{\boldsymbol{\epsilon}}_t^\alpha), \end{aligned} \quad (7.146)$$

where the normal and tangential components of the respective second-order tensors are defined as in Equation (7.145), and

$$\mathbf{C}_{nn}^\alpha := \mathcal{N} : \mathbf{C}^\alpha : \mathcal{N}, \quad \mathbf{C}_{nt}^\alpha := \mathcal{N} : \mathbf{C}^\alpha : \mathcal{T}, \quad \mathbf{C}_{tn}^\alpha := \mathcal{T} : \mathbf{C}^\alpha : \mathcal{N}, \quad \mathbf{C}_{tt}^\alpha := \mathcal{T} : \mathbf{C}^\alpha : \mathcal{T}. \quad (7.147)$$

Remark 148. Note that Equation (7.146) is formally completely identical with the stress-calculation as it would be performed in the \mathcal{B} -system based on the subvectors $\boldsymbol{\epsilon}_n^\alpha$ and $\boldsymbol{\epsilon}_t^\alpha$ and subblocks of the stiffness-tensor defined in (7.105). The interpretation of the various quantities is quite different though, as, due to the use of the projection operators, these quantities are now not, as in [64], elements of lower-dimensional subspaces, but instead elements of the original full-dimensional space (but which do of course each actually lie in a lower-dimensional subspace).

While this also induces some technical difficulties, this representation has the advantage of on the one hand not requiring a basis transformation and on the other hand - as seen from the definitions of \mathcal{N} and \mathcal{T} - being constructed based solely on the normal vector \mathbf{n} . \diamond

Due to the strong formal similarities, it is possible to closely mimick the calculations in [64] within this new formulation. The only technical difficulty is that, as in particular \mathbf{C}_{nn}^α now has a non-trivial kernel even in Voigt-notation as it is based upon a pre- and postmultiplication by a (non-trivial) projector, one cannot rely on actual inverses for systems involving \mathbf{C}_{nn}^α . One can, as in Tschukin, replace the (non-existent) inverses with respect to the normal subspace by Moore-Penrose pseudo-inverses though⁸⁰. At least in theory, the lack of actual invertibility does then not cause any significant issues. In particular, using the corresponding pseudo-inverse $\bar{\mathcal{S}}_{nn}^{12}$ of

$$\bar{\mathcal{C}}_{nn}^{12} := h^1(\phi)\mathcal{N} : \mathcal{C}^1 : \mathcal{N} + h^2(\phi)\mathcal{N} : \mathcal{C}^2 : \mathcal{N} = \mathcal{N} : \bar{\mathcal{C}}^{12}(\phi) : \mathcal{N}$$

with the anti-arithmetic average $\bar{\mathcal{C}}^{12}(\phi) = h^1(\phi)\mathcal{C}^1 + h^2(\phi)\mathcal{C}^2$ as in Equation (7.90), he derives the expression⁸¹

$$f_{el} = \frac{1}{2} \left(\boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) h^1(\phi) + \boldsymbol{\Sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) h^2(\phi) - (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) : \bar{\mathcal{S}}_{nn}^{12} : (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) h^1(\phi) h^2(\phi) \right) \quad (7.148)$$

⁸⁰Constructed such as to act as the corresponding inverses when restricted to the normal subspace only, these satisfy in particular $(\cdot)^\dagger : (\cdot) = \mathcal{N}$ as well as $(\cdot)^\dagger : \mathcal{N} = \mathcal{N} : (\cdot)^\dagger = (\cdot)^\dagger$, $(\cdot)^\dagger : \mathcal{T} = \mathcal{T} : (\cdot)^\dagger = \mathbf{0}$, see [74] and the discussion below.

⁸¹Note that, while Tschukin’s starting point is - similarly to the work in [64] - an expression for the elastic free energy in terms of the “homogeneous” variables $\boldsymbol{\sigma}_n$ and $\boldsymbol{\epsilon}_t$, he then proceeds a somewhat different line as a preparation to a quite different generalization to the multiphase case from the one in [62]. This will further be discussed in Section 7.2.4.

for the elastic free energy density, where $\boldsymbol{\Sigma}^\alpha = \boldsymbol{\mathcal{C}}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)$, $\alpha = 1, 2$, corresponds to the stress obtained under the Voigt-Taylor assumption $\boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon}^{\text{82}}$. The effective stress is then, in accordance with Equation (7.70), determined by differentiating f_{el} with respect to the average strain $\boldsymbol{\epsilon}$ to be given by

$$\boldsymbol{\sigma} = \boldsymbol{\Sigma}^1 h^1(\phi) + \boldsymbol{\Sigma}^2 h^2(\phi) - h^1(\phi) h^2(\phi) (\boldsymbol{\mathcal{C}}^1 - \boldsymbol{\mathcal{C}}^2) : \bar{\boldsymbol{\mathcal{S}}}_{nn}^{12} : (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2). \quad (7.149)$$

Even though the model in [74] starts from a similar argument as the one in [64], as one might guess by comparing equations (7.149) and (7.91), it is in fact much easier to relate to the jump-vector based formulation in Section 7.2.2 than to the one in [64]. The key observation to this is contained in the following lemma:

Lemma 9. *Let $\boldsymbol{\mathcal{C}}$ be a symmetric (in the same sense as the standard stiffness-tensor with the major and minor subsymmetries) fourth-order such that $\boldsymbol{n} \cdot \boldsymbol{\mathcal{C}} \cdot \boldsymbol{n}$ is invertible. Then the pseudo-inverse $\boldsymbol{\mathcal{S}}_{nn} := (\mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N})^\dagger$ satisfies*

$$\boldsymbol{\mathcal{S}}_{nn} : \boldsymbol{\Sigma} = (\boldsymbol{a} \otimes \boldsymbol{n})_S \quad \text{where} \quad \boldsymbol{a} = (\boldsymbol{n} \cdot \boldsymbol{\mathcal{C}} \cdot \boldsymbol{n})^{-1} \cdot (\boldsymbol{\Sigma} \cdot \boldsymbol{n}) \quad (7.150)$$

for any symmetric second-order tensor $\boldsymbol{\Sigma}$.

Proof. Even though it would be possible to verify this in Voigt-notation by retransforming to the \mathcal{B} -system and using Remark 130 together with the discussions in Remark 140 and 141, it is actually simpler to directly verify that the construction in Equation (7.150) satisfies all the required properties.

Based on the definitions of the projectors in Equation (7.144), it is easy to verify that $\boldsymbol{\mathcal{S}}_{nn}$ satisfies $\mathcal{T} : \boldsymbol{\mathcal{S}}_{nn} = \boldsymbol{\mathcal{S}}_{nn} : \mathcal{T} = \mathbf{0}$ and $\mathcal{N} : \boldsymbol{\mathcal{S}}_{nn} = \boldsymbol{\mathcal{S}}_{nn} : \mathcal{N} = \boldsymbol{\mathcal{S}}_{nn}$. With the definition in the left-most equality in Equation (7.150), the equalities for $\mathcal{T} : \boldsymbol{\mathcal{S}}_{nn}$ and $\mathcal{N} : \boldsymbol{\mathcal{S}}_{nn}$ are in fact independent of the particular choice of \boldsymbol{a} but simply rely on the fact that $(\boldsymbol{a} \otimes \boldsymbol{n})_S$ is an element of the normal subspace⁸³ and thus satisfies $\mathcal{T} : (\boldsymbol{a} \otimes \boldsymbol{n})_S = \mathbf{0}$ and $\mathcal{N} : (\boldsymbol{a} \otimes \boldsymbol{n})_S = (\boldsymbol{a} \otimes \boldsymbol{n})_S$. Applying \mathcal{T} resp. \mathcal{N} to $\boldsymbol{\mathcal{S}}_{nn} : \boldsymbol{\Sigma} = (\boldsymbol{a} \otimes \boldsymbol{n})_S$ then leads to the desired conclusion.

The corresponding equalities with \mathcal{T} and \mathcal{N} applied from the right in turn depend upon $\boldsymbol{\Sigma}$ entering the formula for \boldsymbol{a} in terms of its normal component $\boldsymbol{\Sigma} \cdot \boldsymbol{n}$ only. Since $(\mathcal{T} : \boldsymbol{\Sigma}) \cdot \boldsymbol{n} = \mathbf{0}$ and $(\mathcal{N} : \boldsymbol{\Sigma}) \cdot \boldsymbol{n} = \boldsymbol{\Sigma} \cdot \boldsymbol{n}$ ⁸⁴, the definition of \boldsymbol{a} shows that $\boldsymbol{\mathcal{S}}_{nn} : \mathcal{T} : \boldsymbol{\Sigma} = \mathbf{0}$ in the first case, whereas \boldsymbol{a} remains unchanged in the second case and thus $\boldsymbol{\mathcal{S}}_{nn} : \mathcal{N} : \boldsymbol{\Sigma} = \boldsymbol{\mathcal{S}}_{nn} : \boldsymbol{\Sigma}$.

As $\boldsymbol{\mathcal{S}}_{nn}$ therefore maps into the normal subspace and its kernel contains (at least) the tangential subspace, all that remains to verify is that $\boldsymbol{\mathcal{S}}_{nn}$ acts as an inverse to $\mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N}$ on the normal subspace (of symmetric tensors), i.e. $\boldsymbol{\mathcal{S}}_{nn} : \mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N} : \boldsymbol{\Sigma} = \mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N} : \boldsymbol{\mathcal{S}}_{nn} : \boldsymbol{\Sigma} = \mathcal{N} : \boldsymbol{\Sigma}$ for all symmetric $\boldsymbol{\Sigma}$. Applying $\mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N}$ from the left to the left equality in equation (7.150), it first follows from $\mathcal{N} : (\boldsymbol{a} \otimes \boldsymbol{n})_S$ and the right subsymmetry of $\boldsymbol{\mathcal{C}}$ that

$$(\mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N}) : \boldsymbol{\mathcal{S}}_{nn} : \boldsymbol{\Sigma} = (\mathcal{N} : \boldsymbol{\mathcal{C}} : \mathcal{N}) : (\boldsymbol{a} \otimes \boldsymbol{n})_S = \mathcal{N} : \boldsymbol{\mathcal{C}} : (\boldsymbol{a} \otimes \boldsymbol{n})_S = \mathcal{N} : ((\boldsymbol{\mathcal{C}} \cdot \boldsymbol{n}) \cdot \boldsymbol{a}).$$

⁸²Note that Tschukin also uses the notation \boldsymbol{E} instead of $\boldsymbol{\epsilon}$ for the strain do distinguish it from its representation in Voigt-notation.

⁸³While this has already been used a number of times, it can also be verified directly based on $\boldsymbol{\Phi} = \boldsymbol{n} \otimes \boldsymbol{n}$ and the definition of \mathcal{T} . One has

$$\mathcal{T} : (\boldsymbol{a} \otimes \boldsymbol{n}) = (\boldsymbol{I} - \boldsymbol{n} \otimes \boldsymbol{n}) \cdot (\boldsymbol{a} \otimes \boldsymbol{n}) \cdot (\boldsymbol{I} - \boldsymbol{n} \otimes \boldsymbol{n}) = (\boldsymbol{I} - \boldsymbol{n} \otimes \boldsymbol{n}) \cdot (\boldsymbol{a} \otimes (\boldsymbol{I} \cdot \boldsymbol{n}) - \|\boldsymbol{n}\|^2 \boldsymbol{a} \otimes \boldsymbol{n}) = \mathbf{0}$$

and similarly $\mathcal{T} : (\boldsymbol{n} \otimes \boldsymbol{a}) = \mathbf{0}$. As $\mathcal{N} = \boldsymbol{I} - \mathcal{T}$, $\mathcal{N} : (\boldsymbol{n} \otimes \boldsymbol{a})_S = (\boldsymbol{n} \otimes \boldsymbol{a})_S - \mathcal{T} : (\boldsymbol{n} \otimes \boldsymbol{a})_S = (\boldsymbol{n} \otimes \boldsymbol{a})_S$.

⁸⁴This is intuitively clear again as the former is the projection onto the tangential subspace and the latter the one onto the normal subspace. It can also easily be verified explicitly by applying the projectors to $\boldsymbol{\Sigma}$ and contracting the result with \boldsymbol{n} .

Inserting $((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a})$ into the definition of \mathcal{N} and using the left subsymmetry of \mathcal{C} , one further has

$$\begin{aligned} \mathcal{N} : ((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) &= (\mathbf{I} \square \Phi + \Phi \square \mathbf{I} - \Phi \square \Phi) : ((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) \\ &= \left(((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) \cdot \mathbf{n} \right) \otimes \mathbf{n} + \mathbf{n} \otimes \left(\mathbf{n} \cdot ((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) \right) - \left(\mathbf{n} \cdot ((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) \cdot \mathbf{n} \right) \mathbf{n} \otimes \mathbf{n} \\ &= \left((\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a} \right) \otimes \mathbf{n} + \mathbf{n} \otimes \left((\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a} \right) - \left((\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a} \cdot \mathbf{n} \right) \mathbf{n} \otimes \mathbf{n}. \end{aligned}$$

Together with the definition of \mathbf{a} above, the prefactors $\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n}$ cancel, leaving

$$\mathcal{N} : ((\mathcal{C} \cdot \mathbf{n}) \cdot \mathbf{a}) = (\boldsymbol{\Sigma} \cdot \mathbf{n}) \otimes \mathbf{n} + \mathbf{n} \otimes (\boldsymbol{\Sigma} \cdot \mathbf{n}) - \left((\boldsymbol{\Sigma} \cdot \mathbf{n}) \cdot \mathbf{n} \right) \mathbf{n} \otimes \mathbf{n}.$$

As this is precisely $\mathcal{N} : \boldsymbol{\Sigma}$, it follows that

$$(\mathcal{N} : \mathcal{C} : \mathcal{N}) : \mathcal{S}_{nn} : \boldsymbol{\Sigma} = \mathcal{N} : \boldsymbol{\Sigma}$$

for any (symmetric) $\boldsymbol{\Sigma}$, and thus the result. \square

Remark 149. Note that the (somewhat tedious) restriction to symmetric tensors in the lemma is due to the fact that it is on the one hand fully sufficient in the current setting and on the other hand avoids a slight technical difficulty related to the various subsymmetries of fourth-order tensors. More precisely, as $(\mathbf{a} \otimes \mathbf{n})_S$ is obviously a symmetric tensor, \mathcal{S}_{nn} as defined above maps into the subspace of symmetric second order tensors and thus enjoys the left subsymmetric $(\mathcal{S}_{nn})_{ijkl} = (\mathcal{S}_{nn})_{jikl}$. One could also enforce the right subsymmetry by replacing the definition of \mathbf{a} by $\mathbf{a} := (\mathbf{n} \cdot \mathcal{C} \cdot \mathbf{n})^{-1} \cdot \frac{1}{2} (\boldsymbol{\Sigma} \cdot \mathbf{n} + \boldsymbol{\Sigma}^T \cdot \mathbf{n})$, which, although not really used in the current setting, would not require any major modifications either and would then indeed lead to the required pseudoinverse.

The only somewhat unpleasant issue in the general case is related to the definition of the projectors \mathcal{N} and \mathcal{T} themselves. In fact, even though the projector \mathcal{N} in the form of Equation (7.144) does **not** possess either the left nor right subsymmetry⁸⁵, it is well suited to the setting here⁸⁶ and avoids introducing unnecessary symmetrization operations. \diamond

Based on Lemma 9, it is now easy to see that, within the two-phase setting, the model by [74] is equivalent to the jump-condition based one in Section 7.2.2 and by the previous discussion thus also (except for the missing contribution due to $\frac{\partial f_{el}}{\partial \nabla \phi}$ in [64]) with the one in [64].

Firstly, using this lemma for replacing the term $\mathcal{S}_{nn}^{12}(\phi)$ in Equation (7.149) with $(\mathbf{a} \otimes \mathbf{n})_S$ with \mathbf{a} defined by $\tilde{\mathcal{C}}_{nn}^{12}(\phi) \cdot \mathbf{a} = (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) \cdot \mathbf{n}$ with the jump-vector \mathbf{a} coinciding with the one in Equation (7.89) and making use of the \mathcal{C}^α , the expression for the effective stress in [74] clearly reduces to the one from Equation (7.91) in Section 7.2.2. Secondly, the same replacement in the expression for the elastic free energy density reduces Equation (7.148) to

$$f_{el} = \frac{1}{2} \left(\boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) h^1(\phi) + \boldsymbol{\Sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) h^2(\phi) - (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) : (\mathbf{a} \otimes \mathbf{n})_S h^1(\phi) h^2(\phi) \right). \quad (7.151)$$

For comparison with the jump-condition based approach, note that, by construction, the elastic free energy density from Equation (7.96) in Section 7.2.2 can be rewritten as

$$\begin{aligned} f_{el} &= \frac{1}{2} \left(\boldsymbol{\sigma}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) h^1(\phi) + \boldsymbol{\sigma}^2 : (\boldsymbol{\epsilon}^2 - \tilde{\boldsymbol{\epsilon}}^2) h^2(\phi) \right) \\ &= \frac{1}{2} \left(\boldsymbol{\sigma}^1 : \left(\boldsymbol{\epsilon} - h^2(\phi) (\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^1 \right) h^1(\phi) + \boldsymbol{\sigma}^2 : \left(\boldsymbol{\epsilon} + h^1(\phi) (\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^2 \right) h^2(\phi) \right). \end{aligned}$$

⁸⁵This can be seen from $\mathcal{N} : \boldsymbol{\Sigma} = (\boldsymbol{\Sigma} \cdot \mathbf{n}) \otimes \mathbf{n} + \mathbf{n} \otimes (\boldsymbol{\Sigma}^T \cdot \mathbf{n}) - (\mathbf{n} \cdot \boldsymbol{\Sigma} \cdot \mathbf{n}) \mathbf{n} \otimes \mathbf{n}$, which is not in general symmetric unless $\boldsymbol{\Sigma}$ is so, nor depends symmetrically on $\boldsymbol{\Sigma}$.

⁸⁶Both \mathcal{N} and \mathcal{T} are here only applied to quantities with the necessary subsymmetries for this not to matter (either symmetric second-order tensors or fourth-order tensors with the left and right subsymmetries).

As already used in various forms above, the terms in $(\mathbf{a} \otimes \mathbf{n})_S$ cancel as they can be combined to $((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{n})h^1(\phi)h^2(\phi)$ which vanishes due to the common normal stress, i.e. an equivalent expression is given by

$$f_{el} = \frac{1}{2} \left(\boldsymbol{\sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) h^1(\phi) + \boldsymbol{\sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) h^2(\phi) \right).$$

Further inserting

$$\begin{aligned} \boldsymbol{\sigma}^1 &= \mathcal{C}^1 : \left(\boldsymbol{\epsilon} - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^1 \right) = \boldsymbol{\Sigma}^1 - h^2(\phi) \mathcal{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S, \\ \boldsymbol{\sigma}^2 &= \mathcal{C}^2 : \left(\boldsymbol{\epsilon} + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^2 \right) = \boldsymbol{\Sigma}^2 + h^1(\phi) \mathcal{C}^2 : (\mathbf{a} \otimes \mathbf{n})_S \end{aligned} \quad (7.152)$$

then leads to

$$f_{el} = \frac{1}{2} \left(\boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) h^1(\phi) + \boldsymbol{\Sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) h^2(\phi) \right) - (\mathbf{a} \otimes \mathbf{n})_S : \left(\mathcal{C}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - \mathcal{C}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) \right) h^1(\phi) h^2(\phi),$$

which is precisely Equation (7.151), thus showing that both free energy densities also coincide⁸⁷.

Based on the equality of f_{el} , it is clear that the contributions $\frac{\partial f_{el}}{\partial \phi}$ and $\frac{\partial f_{el}}{\partial \nabla \phi} = \frac{\partial f_{el}}{\partial \mathbf{n}} \cdot \frac{\partial \mathbf{n}}{\partial \nabla \phi}$ should also coincide. Tschukin specifies the expression for $\frac{\partial f_{el}}{\partial \phi}$ as⁸⁸

$$\frac{\partial f_{el}}{\partial \phi} = \frac{1}{2} \left(\left(\boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1) - \Delta \tilde{\boldsymbol{\epsilon}}^{12} : \mathcal{C}^1 : \Delta \tilde{\boldsymbol{\epsilon}}^{12} (h^2)^2 \right) - \left(\boldsymbol{\Sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) - \Delta \tilde{\boldsymbol{\epsilon}}^{21} : \mathcal{C}^2 : \Delta \tilde{\boldsymbol{\epsilon}}^{21} (h^1)^2 \right) \right) \frac{\partial h^1}{\partial \phi} \quad (7.153)$$

where $\Delta \tilde{\boldsymbol{\epsilon}}^{12}$ and $\Delta \tilde{\boldsymbol{\epsilon}}^{21}$ are specified by

$$\Delta \tilde{\boldsymbol{\epsilon}}^{12} = \bar{\mathcal{S}}_{nn}^{12} : (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) \quad \text{and} \quad \Delta \tilde{\boldsymbol{\epsilon}}^{21} = \bar{\mathcal{S}}_{nn}^{21} : (\boldsymbol{\Sigma}^2 - \boldsymbol{\Sigma}^1).$$

From Lemma 9, and $\bar{\mathcal{S}}_{nn}^{12} = \bar{\mathcal{S}}_{nn}^{21}$, it is easily seen that these quantities in fact correspond to $\Delta \tilde{\boldsymbol{\epsilon}}^{12} = (\mathbf{a} \otimes \mathbf{n})_S$ and $\Delta \tilde{\boldsymbol{\epsilon}}^{21} = -\Delta \tilde{\boldsymbol{\epsilon}}^{12} = -(\mathbf{a} \otimes \mathbf{n})_S$.

In addition, using $\boldsymbol{\epsilon} = \boldsymbol{\epsilon}^1 + h^2(\mathbf{a} \otimes \mathbf{n})_S$ and thus $\boldsymbol{\Sigma}^1 = \boldsymbol{\sigma}^1 + h^2 \mathcal{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S$, one can rewrite the first term $\boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^1)$ as

$$\begin{aligned} & \boldsymbol{\Sigma}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) + h^2 \boldsymbol{\Sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S \\ &= \boldsymbol{\sigma}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) + h^2 (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) : \mathcal{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S + h^2 \boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S + (h^2)^2 (\mathbf{a} \otimes \mathbf{n})_S : \mathcal{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S \\ &= \boldsymbol{\sigma}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) + 2h^2 \boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S + (h^2)^2 (\mathbf{a} \otimes \mathbf{n})_S : \mathcal{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S \end{aligned}$$

and similarly, with $\boldsymbol{\epsilon} = \boldsymbol{\epsilon}^2 - h^1(\mathbf{a} \otimes \mathbf{n})_S$ and thus $\boldsymbol{\Sigma}^2 = \boldsymbol{\sigma}^2 - h^1 \mathcal{C}^2 : (\mathbf{a} \otimes \mathbf{n})_S$,

$$\boldsymbol{\Sigma}^2 : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^2) = \boldsymbol{\sigma}^2 : (\boldsymbol{\epsilon}^2 - \tilde{\boldsymbol{\epsilon}}^2) - 2h^1 \boldsymbol{\sigma}^2 : (\mathbf{a} \otimes \mathbf{n})_S + (h^1)^2 (\mathbf{a} \otimes \mathbf{n})_S : \mathcal{C}^2 : (\mathbf{a} \otimes \mathbf{n})_S.$$

Combining these observations in the expression (7.153), the quadratic terms in $\Delta \tilde{\boldsymbol{\epsilon}}^{12}$ and $\Delta \tilde{\boldsymbol{\epsilon}}^{21}$ cancel, leaving

$$\frac{\partial f_{el}}{\partial \phi} = \frac{1}{2} \left(\left(\boldsymbol{\sigma}^1 : (\boldsymbol{\epsilon}^1 - \tilde{\boldsymbol{\epsilon}}^1) + 2h^2 \boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S \right) - \left(\boldsymbol{\sigma}^2 : (\boldsymbol{\epsilon}^2 - \tilde{\boldsymbol{\epsilon}}^2) - 2h^1 \boldsymbol{\sigma}^2 : (\mathbf{a} \otimes \mathbf{n})_S \right) \right) \frac{\partial h^1}{\partial \phi}.$$

⁸⁷Note that the equality of the stresses could also have been derived from the equality of the free energy densities, as the effective stress in both cases satisfies $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$. Nevertheless, the direct verification of the equality of the stresses is much simpler.

⁸⁸Note that the expressions are slightly adopted as compared to the ones in [74] since, in contrast to the discussion in Section 7.2.2, Tschukin does not explicitly use a reduced formulation $\phi^2 = 1 - \phi^1$ but instead works with $h^1(\phi^1)$ and $h^2(\phi^2)$ with ϕ^1 and ϕ^2 considered as independent. Once the Lagrange-multiplier for the sum-constraint is taken into account, this reduces, up to the usual factor 2, to the same.

While the terms $\frac{1}{2}\boldsymbol{\sigma}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ correspond to the phase-specific contributions f_{el}^α to f_{el} , the remaining terms in the jump vector \mathbf{a} can, making use of the equality of the normals stresses, be summarized to

$$h^2\boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S + h^1\boldsymbol{\sigma}^2 : (\mathbf{a} \otimes \mathbf{n})_S = (h^2\boldsymbol{\sigma}^1 \cdot \mathbf{n} + h^1\boldsymbol{\sigma}^2 \cdot \mathbf{n}) \cdot \mathbf{a} = (h^2 + h^1)\boldsymbol{\sigma} \cdot \mathbf{a} = (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{a},$$

showing that the expression for $\frac{\partial f_{el}}{\partial \phi}$ coincides with the one in Equation (7.97) (and therefore of course also with the other alternative expressions listed in Section 7.2.2).

Tschukin also specifies the (auxiliary but decisive) quantity $\frac{\partial f_{el}}{\partial n_i}$ in terms of $\Delta\bar{\boldsymbol{\epsilon}}^{12}$ as

$$\frac{\partial f_{el}}{\partial n_i} = -(\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) : (\mathcal{D}_i - \bar{\boldsymbol{S}}_{nn}^{12} : \mathcal{D}_i : \bar{\boldsymbol{C}}^{12}) : \Delta\bar{\boldsymbol{\epsilon}}^{12} h^1(\phi) h^2(\phi) \quad (7.154)$$

where $\mathcal{D}_i = \frac{\partial}{\partial n_i} \mathcal{N}$. Noting first that $\bar{\boldsymbol{S}}_{nn}^{12} : (\boldsymbol{\Sigma}^1 - \boldsymbol{\Sigma}^2) = \Delta\bar{\boldsymbol{\epsilon}}^{12}$, this can, together with the major symmetry of \mathcal{D}_i , be rewritten as⁸⁹

$$\frac{\partial f_{el}}{\partial n_i} = (\boldsymbol{\Sigma}^2 - \boldsymbol{\Sigma}^1 - \bar{\boldsymbol{C}}^{12} : \Delta\bar{\boldsymbol{\epsilon}}^{12}) : \mathcal{D}_i : \Delta\bar{\boldsymbol{\epsilon}}^{12} h^1(\phi) h^2(\phi).$$

Replacing $\Delta\bar{\boldsymbol{\epsilon}}^{12}$ as above with $(\mathbf{a} \otimes \mathbf{n})_S$ and with $\bar{\boldsymbol{C}}^{12}(\phi) = h^2(\phi)\boldsymbol{C}^1 + h^1(\phi)\boldsymbol{C}^2$, the first term can be rewritten as

$$\boldsymbol{\Sigma}^2 - \boldsymbol{\Sigma}^1 + \bar{\boldsymbol{C}}^{12} : (\mathbf{a} \otimes \mathbf{n})_S = (\boldsymbol{\Sigma}^2 + h^1\boldsymbol{C}^2 : (\mathbf{a} \otimes \mathbf{n})_S) - (\boldsymbol{\Sigma}^1 - h^2\boldsymbol{C}^1 : (\mathbf{a} \otimes \mathbf{n})_S) = \boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1,$$

where the final equality follows from Equation (7.152). This further simplifies $\frac{\partial f_{el}}{\partial n_i}$ to $\frac{\partial f_{el}}{\partial n_i} = (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) : \mathcal{D}_i (\mathbf{a} \otimes \mathbf{n})_S$. As a final step for relating the expression (7.154) with the one previously derived in Equation (7.99) for the jump-vector based approach, it remains to eliminate the derivatives \mathcal{D}_i of the projector \mathcal{N} . Firstly, by a simple application of the chain-rule,

$$\mathcal{D}_i : (\mathbf{a} \otimes \mathbf{n})_S = \left(\frac{\partial}{\partial n_i} \mathcal{N} \right) : (\mathbf{a} \otimes \mathbf{n})_S = \frac{\partial}{\partial n_i} (\mathcal{N} (\mathbf{a} \otimes \mathbf{n})_S) - \mathcal{N} : \frac{\partial}{\partial n_i} (\mathbf{a} \otimes \mathbf{n})_S$$

Secondly, $(\mathbf{a} \otimes \mathbf{n})_S$ being an element of the normal subspace, $\mathcal{N} : (\mathbf{a} \otimes \mathbf{n})_S = (\mathbf{a} \otimes \mathbf{n})_S$, and thus

$$\mathcal{D}_i : (\mathbf{a} \otimes \mathbf{n})_S = \frac{\partial}{\partial n_i} (\mathbf{a} \otimes \mathbf{n})_S - \mathcal{N} : \frac{\partial}{\partial n_i} (\mathbf{a} \otimes \mathbf{n})_S = \mathcal{T} : \frac{\partial}{\partial n_i} (\mathbf{a} \otimes \mathbf{n})_S.$$

Even though both \mathbf{a} and \mathbf{n} depend on \mathbf{n} , the presence of the tangential projector then ensures that the contribution $\left(\frac{\partial \mathbf{a}}{\partial n_i} \otimes \mathbf{n} \right)_S$ due to the derivative of \mathbf{a} drops out as this is (regardless of the entries of the first vector in this dyad), again an element of the normal space. This further simplifies the expression for $\mathcal{D}_i : (\mathbf{a} \otimes \mathbf{n})_S$ to

$$\mathcal{D}_i : (\mathbf{a} \otimes \mathbf{n})_S = \mathcal{T} : \left(\mathbf{a} \otimes \frac{\partial \mathbf{n}}{\partial n_i} \right)_S = \mathcal{T} : (\mathbf{a} \otimes \mathbf{e}_i)_S.$$

Reinserting this into the expression for $\frac{\partial f_{el}}{\partial n_i}$, as the last step, it suffices to realize that by the normal continuity of the stresses, $\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1$ is an element of the tangential subspace and thus satisfies $\mathcal{T} : (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) = \boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1$, from which it then follows that

$$\frac{\partial f_{el}}{\partial n_i} = (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) : \mathcal{T} : (\mathbf{a} \otimes \mathbf{e}_i)_S = \left(\mathcal{T} : (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \right) : (\mathbf{a} \otimes \mathbf{e}_i)_S = (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) : (\mathbf{a} \otimes \mathbf{e}_i)_S.$$

⁸⁹This expression corresponds to the non-Voigt version of the alternative expression for $\frac{\partial f_{el}}{\partial n_i}$ provided by equation (4.58) in [74].

Together with the symmetry of the first term, this is the same as $\left((\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}\right) \cdot \mathbf{e}_i$, i.e. the derivative of f_{el} with respect to n_i is simply the i 'th of the vector $(\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}$ and therefore

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = (\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}. \quad (7.155)$$

Contracting this expression with $\frac{\partial \mathbf{n}}{\partial \nabla \phi}$, this obviously leads to the same expression for $\frac{\partial f_{el}}{\partial \nabla \phi}$ as in Equation (7.100).

Remark 150. Note that Tschukin also provides a description of his model in Voigt-representation, where he includes some additional simplifications as compared to the non-Voigt description above. Some of these are in part similar to the manipulations above, whereas others are more similar to the ones performed when discussing the model in [64].

In particular, he factorizes the Voigt-representation $\mathbf{N}_{\boldsymbol{\sigma}}^v$ of the projector onto the normal space as $\mathbf{N}_{\boldsymbol{\sigma}} = \mathbf{A}\mathbf{B}$, where $\mathbf{A} \in \mathbf{R}^{6 \times 3}$ and $\mathbf{B} \in \mathbf{R}^{3 \times 6}$ (as above) is the matrix representing the multiplication of a stress-type tensor in Voigt-notation by the normal vector \mathbf{n} . Recalling Remark 140, $\mathbf{N}_{\boldsymbol{\sigma}}^v$ (in the two-dimensional setting) and making use of Equation (7.118), this projector can also be written as

$$\mathbf{N}_{\boldsymbol{\sigma}}^v := \mathbf{M}_{\boldsymbol{\epsilon}}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \mathbf{M}_{\boldsymbol{\sigma}} = \mathbf{M}_{\boldsymbol{\epsilon}}^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix} \mathbf{Q}\mathbf{B},$$

from which it follows that \mathbf{A} corresponds precisely to the product of the first three matrices⁹⁰. He then makes use of the idempotence of $\mathbf{N}_{\boldsymbol{\sigma}}^v$ and the “good guess” for the required pseudo-inverse $\mathbf{S}_{\mathbf{nn}}^{12}$ in Voigt-nation for finally eliminating \mathbf{A} completely, reducing e.g. his stress-calculation to a form much more similar to the one in Equation (7.124).

A particular consequence of this is that the matrices \mathbf{D}_i introduced in this Voigt-setting should not be confused with the derivatives of the projector above as these are actually defined as the derivatives of the matrix \mathbf{B} with respect to n_i . As this is just a matrix-form of encoding the multiplication by \mathbf{n} , it is obvious that $\frac{\partial \mathbf{B}}{\partial n_i}$ is simply a matrix-form of encoding a multiplication by the unit vector \mathbf{e}_i and that his expression for the derivative in $\frac{\partial f_{el}}{\partial n_i}$ in equation 4.58 thus in fact just corresponds to the i 'th component of $(\boldsymbol{\sigma}^2 - \boldsymbol{\sigma}^1) \cdot \mathbf{a}$ in Equation (7.155). \diamond

⁹⁰One major difference being that Tschukin reexpresses this in terms of the normal vector only.

7.2.4 Some Partial Extension to the Multiphase Case

Given that the jump-condition based model delivers very satisfactory results in the two-phase case, it is natural to try to extend it to the more general multiphase setting. Unfortunately, in contrast to e.g. the simpler Voigt-Talor or Reuss-Sachs approach, the extension of the more elaborate model described in the previous section faces a very fundamental difficulty which is not a priori related to the phasefield modeling approach itself. Rather, it is related to the description of the mechanical behavior at irregular points such as corner points or, within the present context, points where multiple interfaces meet, since it is well known that - unlike for flat or smoothly curved interfaces - such points (and/or edges in the three-dimensional setting) are generally associated with singularities in the behavior of the solution.

In relation with the jump condition approach from the previous section and the following discussion, it is instructive to illustrate this difficulty in a more concrete form in the example in Figure 7.1, which, for notational simplicity, is based on the analogous approach for the linear heat-conduction problem

$$\begin{cases} \nabla \cdot \mathbf{q} = f, & \mathbf{q} = -\kappa \nabla T & \text{in } \Omega \\ T = 0 & & \text{on } \partial\Omega \end{cases}$$

In this case, the continuity of the traces (and their tangential derivatives) leads to the jump condition $\nabla T^\beta = \nabla T^\alpha + a^{\alpha\beta} \mathbf{n}^{\alpha\beta}$, whereas for f regular (e.g. in $L^2(\Omega)$), the continuity of the normal stress-components needs to be replaced by $\llbracket \mathbf{q}^\alpha - \mathbf{q}^\beta \rrbracket \cdot \mathbf{n}^{\alpha\beta} = 0$.

Simply **assuming** that these jump conditions between the two-phase interfaces continue to hold as one approaches the triple point (i.e. there is a well-defined limit, regardless along which interface this point is approached), the jump condition on the 2-3-interface and the combination of the two jump-conditions for ∇T^1 at the 1-2- and 1-3-interface lead to

$$\nabla T^3 = \nabla T^2 + a^{23} \mathbf{e}_x \quad \text{and} \quad \nabla T^2 + a^{21} \mathbf{e}_y = \nabla T^1 = \nabla T^3 + a^{31} \mathbf{e}_y,$$

implying that $a^{21} \mathbf{e}_y = a^{23} \mathbf{e}_x + a^{31} \mathbf{e}_y$. By linear independence, this shows that $a^{23} = 0$, i.e. (as expected from the continuity of the tangential derivatives along the 1-2- and 1-3-interface) $\frac{\partial T^1}{\partial x} = \frac{\partial T^2}{\partial x} = \frac{\partial T^3}{\partial x}$, and $a^{21} = a^{31}$, i.e. (as expected from the continuity of the tangential derivative along the 2-3-interface) $\frac{\partial T^2}{\partial y} = \frac{\partial T^1}{\partial y} - a^{21} \mathbf{e}_y = \frac{\partial T^1}{\partial y} - a^{31} \mathbf{e}_y = \frac{\partial T^3}{\partial y}$. Inserting these equalities into the jump conditions on the normal fluxes further leads to the conclusion that

$$\kappa^2 \frac{\partial T^2}{\partial x} = \kappa^3 \frac{\partial T^3}{\partial x} \quad \text{and} \quad \kappa^2 \frac{\partial T^2}{\partial y} = \kappa^2 \left(\frac{\partial T^1}{\partial y} - a^{21} \right) = \kappa^3 \left(\frac{\partial T^1}{\partial y} - a^{31} \right) = \kappa^3 \frac{\partial T^3}{\partial y}.$$

As $\nabla T^2 = \nabla T^3$ by the above, these equality is possible iff either $\kappa^2 = \kappa^3$ (i.e. there is no “real” triple-junction) or both gradients happen to vanish at this point. In the latter case, one necessarily also has $\frac{\partial T^1}{\partial x} = 0$ by continuity of the tangential derivatives and $\frac{\partial T^1}{\partial y} = 0$ due to the equality of the normal fluxes. In summary, the jump conditions between all three phases in this case can only be satisfied when $\kappa^2 = \kappa^3$ or when all three gradients happen to vanish at this point⁹¹.

Remark 151. Note that, from a phasefield perspective, the issue at internal corners does not a priori seem to impose any particular difficulties from a pure modeling point of view as such corners are usually “smoothed out” due to the diffuse interface profile. This seeming simplicity is somewhat misleading in the same sense that, even though e.g. a standard (sharp-interface) FEM-simulation does not in principle require any special treatment of corners points, the accuracy of the resulting approximation will in general be significantly reduced as compared to a “smooth”

⁹¹In contrast, when either $\kappa^2 = \kappa^1$ or $\kappa^3 = \kappa^1$, one would still expect a singular behavior due to the presence of a corner between the two “effective” material regions.

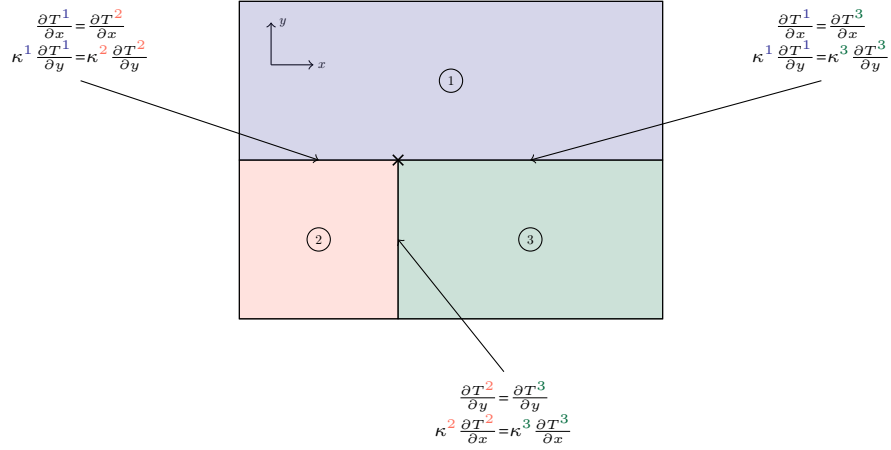


Figure 7.1: Illustration of the contradictory nature of the jump-conditions at a triple point.

problem. Understanding this difficulty is, at least at a qualitative level, quite intuitive, as approximation a solution which has a highly complex behavior within a relatively small region is of course much harder than approximating one which is highly regular.

In contrast to the case of internal corners, which can in principle be tackled using the jump-condition based model above, the case of intersecting interfaces also poses a “visible” difficulty in the sense that the modeling of such multiple junctions within the phasefield context requires using more than two coexistent phases, and thus requires some extension of the previous model(s) to the multiphase case. \diamond

A very similar (though algebraically more challenging) difficulty arises when trying to extend the jump-condition based model from the previous section to the multiphase case. A straightforward generalization of the jump relations on the symmetrized displacement gradient from Equation (7.87) would consist in imposing

$$[[\mathbf{F}]]^{\alpha\beta} = \mathbf{F}^\beta - \mathbf{F}^\alpha = \mathbf{a}^{\alpha\beta} \otimes \mathbf{N}^{\alpha\beta} \quad \text{resp.} \quad [[\boldsymbol{\epsilon}]]^{\alpha\beta} = \boldsymbol{\epsilon}^\beta - \boldsymbol{\epsilon}^\alpha = (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S \quad (7.156)$$

for all $\alpha \neq \beta$ and then trying to determine the $\mathbf{a}^{\alpha\beta}$, $\alpha \neq \beta$ from the corresponding continuity conditions

$$[[\boldsymbol{\sigma}]]^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta} = (\boldsymbol{\sigma}^\beta - \boldsymbol{\sigma}^\alpha) \cdot \mathbf{n}^{\alpha\beta} = \mathbf{0}, \quad \alpha \neq \beta \quad (7.157)$$

for the normal stresses.

As above, such a straightforward generalization is unfortunately doomed to failure (except for some special cases). In fact, even though Equation (7.156) would a priori introduces $N(N-1)$ unknown jump vectors $(\mathbf{a}^{\alpha\beta})_{1 \leq \alpha \neq \beta \leq N}$ together with an equal number of continuity conditions, these conditions are not all independent. Firstly, using $\mathbf{n}^{\beta\alpha} = -\mathbf{n}^{\alpha\beta}$ and thus from

$$\mathbf{a}^{\beta\alpha} \otimes \mathbf{n}^{\beta\alpha} = -\mathbf{a}^{\beta\alpha} \otimes \mathbf{n}^{\alpha\beta} = [[\mathbf{F}]]^{\beta\alpha} = -[[\mathbf{F}]]^{\alpha\beta} = -\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta}$$

that $\mathbf{a}^{\beta\alpha} = \mathbf{a}^{\alpha\beta}$, the number of unknowns can be reduced to the $\frac{N(N-1)}{2}$ unknown jump vectors $(\mathbf{a}^{\alpha\beta})_{1 \leq \alpha < \beta \leq N}$. At the same time, as $[[\boldsymbol{\sigma}]]^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta} = \mathbf{0}$ directly implies $[[\boldsymbol{\sigma}]]^{\beta\alpha} \cdot \mathbf{n}^{\beta\alpha}$, one can simultaneously reduce the number of jump conditions to be satisfied in Equation (7.157) to those for $1 \leq \alpha < \beta \leq N$.

The number of truly independent entries for the jump vectors can be significantly lower though as the algebraic equality

$$[[\mathbf{F}]]^{\alpha\beta} = \mathbf{F}^\beta - \mathbf{F}^\alpha = \mathbf{F}^\beta - \mathbf{F}^\delta + \mathbf{F}^\delta - \mathbf{F}^\alpha = [[\mathbf{F}]]^{\delta\beta} - [[\mathbf{F}]]^{\delta\alpha} = [[\mathbf{F}]]^{\delta\beta} - [[\mathbf{F}]]^{\delta\alpha}$$

imposes the additional compatibility conditions

$$\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta} = \mathbf{a}^{\delta\beta} \otimes \mathbf{n}^{\delta\beta} - \mathbf{a}^{\delta\alpha} \otimes \mathbf{n}^{\delta\alpha} \quad \forall \delta \neq \alpha, \beta \quad (7.158)$$

on the $\mathbf{a}^{\alpha\beta}$. Since $\mathbf{n}^{\alpha\beta}$ is a unit vector, right-multiplying this equation with $\mathbf{n}^{\alpha\beta}$ directly fixes $\mathbf{a}^{\alpha\beta}$ in terms of $\mathbf{a}^{\alpha\delta}$ and $\mathbf{a}^{\beta\delta}$ as⁹²

$$\mathbf{a}^{\alpha\beta} = (\mathbf{n}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}) \mathbf{a}^{\alpha\beta} = (\mathbf{n}^{\alpha\beta} \cdot \mathbf{n}^{\delta\beta}) \mathbf{a}^{\delta\beta} - (\mathbf{n}^{\alpha\beta} \cdot \mathbf{n}^{\delta\alpha}) \mathbf{a}^{\delta\alpha} \quad \forall \delta \neq \alpha.$$

This already implies that **all** jump vectors would immediately follow once they are known with respect to a single phase⁹³.

Even when all normals happen to be parallel (and “testing” the compatibility conditions (7.158) with the common normal vector $\mathbf{n}^{\alpha\beta} = \pm \mathbf{n}$, $\alpha \neq \beta$, therefore provides all relevant information contained in the system) already reduces the number of effectively independent jump vectors to at most $N - 1$ (e.g. the $\mathbf{a}^{1\alpha}$, $\alpha \geq 2$). For every triplet (α, β, δ) of phases for which the normals $\mathbf{n}^{\alpha\beta}$, $\mathbf{n}^{\delta\alpha}$ and $\mathbf{n}^{\delta\beta}$ are not all parallel, Equation (7.158) imposes an additional set of restrictions on the $\mathbf{a}^{\alpha\beta}$ though. In fact, if either $\mathbf{n}^{\delta\alpha}$ or $\mathbf{n}^{\delta\beta}$ are not parallel to $\mathbf{n}^{\alpha\beta}$, there is a tangential vector $\mathbf{t}^{\alpha\beta}$ such that $\mathbf{t}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}$ is zero but at least one of $\mathbf{n}^{\delta\beta} \cdot \mathbf{n}^{\delta\alpha}$ and $\mathbf{t}^{\alpha\beta} \cdot \mathbf{n}^{\delta\beta}$ is non-zero. Applying $\mathbf{t}^{\alpha\beta}$ from the right to the compatibility conditions,

$$\mathbf{0} = (\mathbf{t}^{\alpha\beta} \cdot \mathbf{n}^{\delta\beta}) \mathbf{a}^{\delta\beta} - (\mathbf{t}^{\alpha\beta} \cdot \mathbf{n}^{\delta\alpha}) \mathbf{a}^{\delta\alpha}$$

which allows, besides fixing $\mathbf{a}^{\alpha\beta}$ in terms of the two respective other jump vectors, fixing one jump vector on the right-hand side in terms of the other one (or directly as zero if one of the other normals also happens to be parallel to $\mathbf{n}^{\alpha\beta}$). Having non-parallel “two-phase-normal vectors” - such as is normally the case - thus even further reduces the available degrees of freedom in the ansatz (7.156).

Just as the jumps in the strains are not all independent, neither are of course the jumps in the normal stresses. Nevertheless, this does - similar to the example from Figure 7.1 - not usually entail a sufficient reduction of the number of normal continuity equations to be satisfied, with one notable exception being the case when there is only a single normal involved. In combination with the continuity condition, the resulting system is therefore generally inconsistent since there are not enough degrees of freedom to enforce all equations simultaneously.

Remark 152. Together with the averaging condition $\mathbf{F} = \sum_{\alpha=1}^N \mathbf{F}^{\alpha} h^{\alpha}(\phi)$, it is, similarly to the two-phase setting, possible to reexpress \mathbf{F}^{α} in terms of the average displacement gradient \mathbf{F} and the jumps with respect to the other phases. In fact, making use of $\mathbf{F}^{\beta} = \mathbf{F}^{\alpha} + \llbracket \mathbf{F} \rrbracket^{\alpha\beta}$ and $\sum_{\alpha=1}^N h^{\alpha}(\phi) = 1$ allows rewriting the averaging condition as

$$\mathbf{F} = \mathbf{F}^{\alpha} h^{\alpha}(\phi) + \sum_{\beta \neq \alpha} \mathbf{F}^{\beta} h^{\beta}(\phi) = \mathbf{F}^{\alpha} h^{\alpha}(\phi) + \sum_{\beta \neq \alpha} \llbracket \mathbf{F} \rrbracket^{\alpha\beta} h^{\beta}(\phi) = \mathbf{F} + \sum_{\beta \neq \alpha} \llbracket \mathbf{F} \rrbracket^{\alpha\beta} h^{\beta}(\phi)$$

⁹²Note that, even though this argument is slightly simpler in the non-symmetrized case, the symmetrized form

$$(\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S = (\mathbf{a}^{\delta\beta} \otimes \mathbf{n}^{\delta\beta})_S - (\mathbf{a}^{\delta\alpha} \otimes \mathbf{n}^{\delta\alpha})_S$$

leads to the same conclusion. In fact, a first contraction of the equality $(\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S = \mathbf{T}$ with known \mathbf{T} leads to $\frac{1}{2}(\mathbf{a}^{\alpha\beta} + (\mathbf{a}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}) \mathbf{n}^{\alpha\beta}) = \mathbf{T} \cdot \mathbf{n}^{\alpha\beta}$, which is not yet fully explicit in $\mathbf{a}^{\alpha\beta}$. A second contraction then shows $\mathbf{a}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta} = \mathbf{n}^{\alpha\beta} \cdot \mathbf{T} \cdot \mathbf{n}^{\alpha\beta}$, which, reinserted into the previous equation shows that $\mathbf{a}^{\alpha\beta} = 2\mathbf{T} \cdot \mathbf{n}^{\alpha\beta} - (\mathbf{n}^{\alpha\beta} \cdot \mathbf{T} \cdot \mathbf{n}^{\alpha\beta}) \mathbf{n}^{\alpha\beta}$ and thus fully specifies $\mathbf{a}^{\alpha\beta}$ in terms of \mathbf{T} .

⁹³In fact, assuming for convenience for this to be the first phase, knowing the $\mathbf{a}^{1\alpha}$, $2 \leq \alpha \leq N$ and using the appropriate column in Equation (7.158) with $\delta = 1$ and $\beta = \alpha + 1$ one can obtain all jump vectors $\mathbf{a}^{\alpha\beta}$ for $\alpha = \beta - 1$. Thus knowing $\mathbf{a}^{\alpha(\alpha+1)}$ and $\mathbf{a}^{(\alpha+1)(\alpha+2)}$ for all $\alpha \geq 2$, the $\mathbf{a}^{\alpha(\alpha+2)}$ follow immediately from $\mathbf{a}^{\beta\alpha} = \mathbf{a}^{\alpha\beta}$ and $\mathbf{a}^{\alpha(\alpha+2)} \otimes \mathbf{n}^{\alpha(\alpha+2)} = \mathbf{a}^{(\alpha+1)(\alpha+2)} \otimes \mathbf{n}^{(\alpha+1)(\alpha+2)} - \mathbf{a}^{(\alpha+1)\alpha} \otimes \mathbf{n}^{(\alpha+1)\alpha}$ and, repeating this argument, finally all $\mathbf{a}^{\alpha\beta}$ for $\beta \neq \alpha$.

and thus leads to⁹⁴

$$\mathbf{F}^\alpha = \mathbf{F} - \sum_{\beta \neq \alpha} [\mathbf{F}]^{\alpha\beta} h^\beta(\phi), \forall \alpha, \quad (7.159)$$

respectively, through a completely analogous calculation, the relation

$$\boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon} - \sum_{\beta \neq \alpha} [\boldsymbol{\epsilon}]^{\alpha\beta} h^\beta(\phi), \forall \alpha \quad (7.160)$$

for the $\boldsymbol{\epsilon}^\alpha$. ◇

Remark 153. Note that, taking the weighted average of Equation (7.159) and Equation (7.160) shows that the jumps in addition satisfy

$$\sum_{\alpha=1}^N \sum_{\beta \neq \alpha} [\mathbf{F}]^{\alpha\beta} h^\alpha(\phi) h^\beta(\phi) = \mathbf{0} \quad \text{and} \quad \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} [\boldsymbol{\epsilon}]^{\alpha\beta} h^\alpha(\phi) h^\beta(\phi) = \mathbf{0}. \quad (7.161)$$

◇

A Model Based on a Common Normal ([62], §4.5 [61])

A first way of generalizing the two-phase model from the previous section to the multiphase setting is proposed in [62] and [61] and further successfully applied in e.g. [66], [4], [5]. It is based upon generalizing the representation of the normal and tangential stresses in equations (7.111) and (7.112) in terms of a basis-transformation in equation through the use of a common ‘‘average’’ normal vector \mathbf{n} . Using this normal vector, it assumes the equality of all tangential strains with the (given) average one $\boldsymbol{\epsilon}_t$ and then enforces the equality of all normal phase-specific stresses with the effective (a priori unknown) normal stress $\boldsymbol{\sigma}_n$. Using, instead of the base-transformation as in [62] the projector-based formalism in [74] and setting $(\cdot)_n := \mathcal{N} : (\cdot)$, $(\cdot)_t = \mathcal{T} : (\cdot)$, the underlying assumption is thus that $\boldsymbol{\epsilon}_n^\alpha = \boldsymbol{\epsilon}_t \quad \forall \alpha$ and that $\boldsymbol{\sigma}_n^\alpha = \boldsymbol{\sigma}_n \quad \forall \alpha$. This is combined with the usual averaging conditions $\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^\alpha h^\alpha(\phi)$ and $\boldsymbol{\sigma} = \sum_{\alpha=1}^N \boldsymbol{\sigma}^\alpha h^\alpha(\phi)$ on the total strains and stresses, where, due to the conditions above, it is clear that the averaging procedure actually only affects the normal resp. tangential components. Together with the phase-specific material law $\boldsymbol{\sigma}^\alpha = \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ and making use of the idempotence of the projection operators ($\mathcal{N} = \mathcal{N}^2$ and $\mathcal{T} = \mathcal{T}^2$), the continuity condition on the normal components of the stresses may be rewritten as

$$\begin{aligned} \boldsymbol{\sigma}_n^\alpha &= \mathcal{N} : \left(\mathbf{C}^\alpha : (\mathcal{N} : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) + \mathcal{T} : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)) \right) = \mathcal{N} : \left(\mathbf{C}^\alpha : (\mathcal{N} : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) + \mathcal{T} : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)) \right) \\ &= \mathbf{C}_{nn}^\alpha : (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t) \stackrel{!}{=} \boldsymbol{\sigma}_n := \mathcal{N} : \boldsymbol{\sigma} \quad \forall \alpha, \end{aligned}$$

where $\mathbf{C}_{nn}^\alpha = \mathcal{N} : \mathbf{C}^\alpha : \mathcal{N}$, $\mathbf{C}_{nt}^\alpha = \mathcal{N} : \mathbf{C}^\alpha : \mathcal{T}$ and use was made of $\boldsymbol{\epsilon}_t^\alpha = \boldsymbol{\epsilon}_t$. Using, as in [74], an appropriate pseudo-inverse $\mathbf{S}_{nn}^\alpha = (\mathbf{C}_{nn}^\alpha)^\dagger$ of \mathbf{C}_{nn}^α , this equation permits reexpressing all unknown normal components of the phase-specific strains in terms of the single common normal stress-component as

$$\boldsymbol{\epsilon}_n^\alpha = \tilde{\boldsymbol{\epsilon}}_n^\alpha + \mathbf{S}_{nn}^\alpha : \left(\boldsymbol{\sigma}_n - \mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha) \right) \quad (7.162)$$

and thus leads, through a weighted average, to the condition

$$\boldsymbol{\epsilon}_n = \sum_{\alpha=1}^N \tilde{\boldsymbol{\epsilon}}_n^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^N \mathbf{S}_{nn}^\alpha : \left(\boldsymbol{\sigma}_n - \mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha) \right) h^\alpha(\phi)$$

⁹⁴Note that if $h^\alpha(\phi) = 0$, \mathbf{F}^α does not appear in the averaging condition in terms of the $(\mathbf{F}^\beta)_{1 \leq \beta \leq N}$ and \mathbf{F} and can therefore in principle be defined arbitrarily regardless of the values of the other \mathbf{F}^β , $\beta \neq \alpha$. In contrast, **given** a set of jump vector $[\mathbf{F}]^{\alpha\beta}$ and \mathbf{F} , the relation (7.159) still has to hold for \mathbf{F}^α , the only difference to the case $h^\alpha(\phi) \neq 0$ being that, given all other \mathbf{F}^β and \mathbf{F} , one is free to choose an arbitrary \mathbf{F}^α for defining the jumps $[\mathbf{F}]^{\alpha\beta} = \mathbf{F}^\beta - \mathbf{F}^\alpha$.

relating the $\boldsymbol{\sigma}_n$ with the given average normal strain $\boldsymbol{\epsilon}_n$. Together with the corresponding pseudo-inverse of $\sum_{\alpha=1}^N \mathbf{S}_{nn}^\alpha h^\alpha(\phi)$, which, for consistency with the description in [62] will be abbreviated as $-\bar{\mathcal{T}}_{nn}$, i.e. $\bar{\mathcal{T}}_{nn} := -\left(\sum_{\alpha=1}^N \mathbf{S}_{nn}^\alpha h^\alpha(\phi)\right)^\dagger$, this allows obtaining the common normal stress component as

$$\boldsymbol{\sigma}_n = -\bar{\mathcal{T}}_{nn} : \left(\boldsymbol{\epsilon}_n - \sum_{\alpha=1}^N \tilde{\boldsymbol{\epsilon}}_n^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^N \left(\mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) \right) h^\alpha(\phi) \right). \quad (7.163)$$

Finally, again using Equation (7.162) with this now known normal stress, one can first obtain the normal phase-specific strain components $\boldsymbol{\epsilon}_n^\alpha$ given by

$$\boldsymbol{\epsilon}_n^\alpha = \tilde{\boldsymbol{\epsilon}}_n^\alpha + \mathbf{S}_{nn}^\alpha : \left[-\bar{\mathcal{T}}_{nn} : \left(\boldsymbol{\epsilon}_n - \sum_{\beta=1}^N \tilde{\boldsymbol{\epsilon}}_n^\beta h^\beta(\phi) + \sum_{\beta=1}^N \left(\mathbf{C}_{nt}^\beta : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\beta) \right) h^\beta(\phi) \right) - \mathbf{C}_{nt}^\alpha : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) \right]$$

and, combined with the already known tangential strains $\boldsymbol{\epsilon}_t^\alpha = \boldsymbol{\epsilon}_t$ and through another averaging using the phase-specific material law, the tangential stress components as

$$\begin{aligned} \boldsymbol{\sigma}_t &= \sum_{\alpha=1}^N \boldsymbol{\sigma}_t^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^N \left(\mathbf{C}_{tn}^\alpha : (\boldsymbol{\epsilon}_n^\alpha - \tilde{\boldsymbol{\epsilon}}_n^\alpha) + \mathbf{C}_{tt}^\alpha : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) \right) h^\alpha(\phi) \\ &= - \sum_{\alpha=1}^N \left(\mathbf{C}_{tn}^\alpha : \mathbf{S}_{nn}^\alpha h^\alpha(\phi) \right) : \bar{\mathcal{T}}_{nn} : \left(\boldsymbol{\epsilon}_n - \sum_{\beta=1}^N \tilde{\boldsymbol{\epsilon}}_n^\beta h^\beta(\phi) \right) \\ &\quad + \sum_{\alpha=1}^N \left(\left(\mathbf{C}_{tt}^\alpha - \mathbf{C}_{tn}^\alpha : \mathbf{S}_{nn}^\alpha \mathbf{C}_{nt}^\alpha \right) - \sum_{\beta=1}^N \left(\mathbf{C}_{tn}^\beta : \mathbf{S}_{nn}^\beta h^\beta(\phi) \right) : \bar{\mathcal{T}}_{nn} : \mathbf{C}_{nt}^\alpha \right) : (\boldsymbol{\epsilon}_t - \tilde{\boldsymbol{\epsilon}}_t^\alpha) h^\alpha(\phi) \end{aligned} \quad (7.164)$$

As in the two-phase case, one can alternatively choose a description in terms of appropriately chosen jump vectors. As $\mathcal{T} : \boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon}_t^\alpha = \boldsymbol{\epsilon}_t$ for all α , $\boldsymbol{\epsilon}^\alpha$ is necessarily of the form

$$\boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon} - (\mathbf{a}^\alpha \otimes \mathbf{n})_S. \quad (7.165)$$

Together with the averaging condition $\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^\alpha h^\alpha(\phi)$, this in addition implies the constraint

$$\sum_{\alpha=1}^N h^\alpha(\phi) \mathbf{a}^\alpha \stackrel{!}{=} \mathbf{0} \quad (7.166)$$

on the jump-vectors \mathbf{a}^α .

In terms of the jump vectors and again using the abbreviation $\boldsymbol{\Sigma}^\alpha = \mathbf{C}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)$, the condition on the equality of the normal stresses can be written as

$$\boldsymbol{\sigma}^\alpha \cdot \mathbf{n} = \mathbf{n} \cdot \left(\mathbf{C}^\alpha : (\boldsymbol{\epsilon} - (\mathbf{a}^\alpha \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^\alpha) \right) = \boldsymbol{\Sigma}^\alpha \cdot \mathbf{n} - (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n}) \cdot \mathbf{a}^\alpha \stackrel{!}{=} \boldsymbol{\sigma} \cdot \mathbf{n}$$

and thus shows that the phase-specific jump vectors \mathbf{a}^α are given by

$$\mathbf{a}^\alpha = (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\sigma}) \cdot \mathbf{n} \right) \quad (7.167)$$

in terms of the common value $\boldsymbol{\sigma} \cdot \mathbf{n}$. Together with the zero-average condition on the jump-vectors, it follows that $\boldsymbol{\sigma} \cdot \mathbf{n}$ has to be such that $\sum_{\alpha=1}^N \mathbf{a} h^\alpha(\phi) = \sum_{\alpha=1}^N \left((\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot \left((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\sigma}) \cdot \mathbf{n} \right) \right) h^\alpha(\phi) \stackrel{!}{=} \mathbf{0}$, i.e.

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \left(\sum_{\alpha=1}^N (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} h^\alpha(\phi) \right)^{-1} \cdot \sum_{\alpha=1}^N (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot (\boldsymbol{\Sigma}^\alpha \cdot \mathbf{n}) h^\alpha(\phi). \quad (7.168)$$

Similarly to before, recombining the now known normal stress with Equation (7.167) shows that the phase-specific jump vectors are given by

$$\mathbf{a}^\alpha = (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \cdot \left[(\boldsymbol{\Sigma}^\alpha \cdot \mathbf{n} - \left(\sum_{\beta=1}^N (\mathbf{n} \cdot \mathbf{C}^\beta \cdot \mathbf{n})^{-1} h^\beta(\phi) \right)^{-1} \cdot \sum_{\beta=1}^N (\mathbf{n} \cdot \mathbf{C}^\beta \cdot \mathbf{n})^{-1} \cdot (\boldsymbol{\Sigma}^\beta \cdot \mathbf{n}) h^\beta(\phi) \right],$$

thus leading, with the given value of $\boldsymbol{\epsilon}$ to the alternative expression

$$\begin{aligned} \boldsymbol{\sigma} &= \sum_{\alpha=1}^N \boldsymbol{\Sigma}^\alpha h^\alpha(\phi) + \sum_{\alpha=1}^N (\mathbf{C}^\alpha \cdot \mathbf{n}) \cdot \mathbf{a}^\alpha h^\alpha(\phi) \\ &= \boldsymbol{\Sigma} + \sum_{\alpha=1}^N (\mathbf{C}^\alpha \cdot \mathbf{n}) \cdot (\mathbf{n} \cdot \mathbf{C}^\alpha \cdot \mathbf{n})^{-1} \\ &\quad \cdot \left[(\boldsymbol{\Sigma}^\alpha \cdot \mathbf{n} - \left(\sum_{\beta=1}^N (\mathbf{n} \cdot \mathbf{C}^\beta \cdot \mathbf{n})^{-1} h^\beta(\phi) \right)^{-1} \cdot \sum_{\beta=1}^N (\mathbf{n} \cdot \mathbf{C}^\beta \cdot \mathbf{n})^{-1} \cdot (\boldsymbol{\Sigma}^\beta \cdot \mathbf{n}) h^\beta(\phi) \right] h^\alpha(\phi) \end{aligned} \quad (7.169)$$

for the effective stress.

Based on the use of the jump-vectors above, the driving-force calculation for this model can, alternatively to the procedure in [62] in the transformed coordinate system, be performed in a manner similar to the two-phase case considered previously. Differentiating the expression $f_{el}(\boldsymbol{\phi}, \boldsymbol{\epsilon}) = \sum_{\beta} f_{el}^\beta h^\beta(\boldsymbol{\phi}) = \sum_{\beta} \frac{1}{2} (\boldsymbol{\epsilon}^\beta(\boldsymbol{\phi}, \boldsymbol{\epsilon}) - \tilde{\boldsymbol{\epsilon}}^\beta) : \mathbf{C}^\beta : (\boldsymbol{\epsilon}^\beta(\boldsymbol{\phi}, \boldsymbol{\epsilon}) - \tilde{\boldsymbol{\epsilon}}^\beta) h^\beta(\boldsymbol{\phi})$ with respect to ϕ^α for a **fixed** normal vector \mathbf{n} , it first follows that

$$\frac{\partial}{\partial \phi^\alpha} \Big|_{\mathbf{n}} f_{el}(\boldsymbol{\phi}, \boldsymbol{\epsilon}) = \sum_{\beta} f_{el}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} + \sum_{\beta} \boldsymbol{\sigma}^\beta : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}} h^\beta.$$

In contrast to the chemical case, one can again not directly make use of the product rule and $\frac{\partial \sum_{\beta} \boldsymbol{\epsilon}^\beta h^\beta}{\partial \phi^\alpha} = \frac{\partial \boldsymbol{\epsilon}}{\partial \phi^\alpha} = \mathbf{0}$ for simplifying the last term as the tangential parts of the $\boldsymbol{\sigma}^\beta$ do not reduce to a common value. Nevertheless, since $\boldsymbol{\epsilon}^\beta = \boldsymbol{\epsilon} - (\mathbf{a}^\beta \otimes \mathbf{n})_S$, it follows that $\frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}} = -(\frac{\partial \mathbf{a}^\beta}{\partial \phi^\alpha} \otimes \mathbf{n})_S$ and thus

$$\sum_{\beta} \boldsymbol{\sigma}^\beta : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}} h^\beta = - \sum_{\beta} (\boldsymbol{\sigma}^\beta \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}^\beta}{\partial \phi^\alpha} h^\beta = -(\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \sum_{\beta} \frac{\partial \mathbf{a}^\beta}{\partial \phi^\alpha} h^\beta$$

since all normal stresses for the different phases are in fact equal. Further making use of the condition $\sum_{\beta} \mathbf{a}^\beta h^\beta = \mathbf{0}$ from Equation (7.166) and using the product rule $\sum_{\beta} \frac{\partial \mathbf{a}^\beta}{\partial \phi^\alpha} h^\beta = \frac{\partial \sum_{\beta} \mathbf{a}^\beta h^\beta}{\partial \phi^\alpha} - \sum_{\beta} \mathbf{a}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$, one can eliminate the derivatives of the \mathbf{a}^β , finally leading to

$$\frac{\partial}{\partial \phi^\alpha} \Big|_{\mathbf{n}} f_{el}(\boldsymbol{\phi}, \boldsymbol{\epsilon}) = \sum_{\beta} f_{el}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} + (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \sum_{\beta} \mathbf{a}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} = \sum_{\beta} (f_{el}^\beta + (\boldsymbol{\sigma}^\beta \cdot \mathbf{n}) \cdot \mathbf{a}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha}. \quad (7.170)$$

Alternatively using $\boldsymbol{\epsilon}^\beta = \boldsymbol{\epsilon} - (\mathbf{a}^\beta \otimes \mathbf{n})_S$, one can also rewrite the second contribution to the driving force as

$$\sum_{\beta=1}^N (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{a}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} = \sum_{\beta=1}^N \boldsymbol{\sigma} : (\mathbf{a}^\beta \otimes \mathbf{n})_S \frac{\partial h^\beta}{\partial \phi^\alpha} = \sum_{\beta=1}^N \boldsymbol{\sigma} : (\boldsymbol{\epsilon} - \boldsymbol{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha} = - \sum_{\beta=1}^N \boldsymbol{\sigma} : \boldsymbol{\epsilon}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha} + (\boldsymbol{\sigma} : \boldsymbol{\epsilon}) \frac{\partial \sum_{\beta=1}^N h^\beta}{\partial \phi^\alpha}.$$

Since the last term vanishes, this leaves

$$\frac{\partial f_{el}}{\partial \phi^\alpha} \Big|_{\mathbf{n}} = \sum_{\beta=1}^N (f_{el}^\beta(\boldsymbol{\epsilon}^\beta) - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha}. \quad (7.171)$$

The remaining contributions to $\frac{\partial f_{el}}{\partial \phi^\alpha}$ and $\frac{\partial f_{el}}{\partial \nabla \phi^\alpha}$ through the dependence on \mathbf{n} follow from a similar reasoning in combination with the definition of \mathbf{n} in terms of ϕ and $\nabla \phi$. Differentiating f_{el} with respect to \mathbf{n} one has

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = \sum_{\beta} \boldsymbol{\sigma}^{\beta} : \frac{\partial \boldsymbol{\epsilon}^{\beta}}{\partial \mathbf{n}} h^{\beta} = - \sum_{\beta} \left((\boldsymbol{\sigma}^{\beta} \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}^{\beta}}{\partial \mathbf{n}} + (\boldsymbol{\sigma}^{\beta} \cdot \mathbf{a}^{\beta}) \right) h^{\beta}.$$

The first contribution can be eliminated by making use of the common normal stress $\boldsymbol{\sigma}^{\beta} \cdot \mathbf{n} = \boldsymbol{\sigma} \cdot \mathbf{n}$ for all β and the condition $\sum_{\beta} \mathbf{a}^{\beta} h^{\beta} = \mathbf{0}$ together with the independence of the h^{β} on \mathbf{n} , thus leaving the simple expression

$$\frac{\partial f_{el}}{\partial \mathbf{n}} = - \sum_{\beta=1}^N (\boldsymbol{\sigma}^{\beta} \cdot \mathbf{a}^{\beta}) h^{\beta}(\phi) \quad (7.172)$$

for the derivative with respect to \mathbf{n} .

Remark 154. As the similarity of the calculations and the resulting expressions for the driving force contributions to the two-phase case shows, this model is able to maintain the most of the pleasant properties of its two-phase analogon in the multiphase case. In particular, its free energy density is defined in the natural manner as the weighted average of the $f^{\alpha} = \frac{1}{2}(\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}) : \mathcal{C}^{\alpha} : (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha})$ of the phase-specific ones with the $\boldsymbol{\epsilon}^{\alpha}$ satisfying the Hadamard jump conditions with respect to every other phase. Furthermore, the normal stresses between all phases are continuous, and the total stress satisfies both $\boldsymbol{\sigma} = \sum_{\alpha} \boldsymbol{\sigma}^{\alpha} h^{\alpha}(\phi)$ and the compatibility condition $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$. In fact, one has $\frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}} = \sum_{\alpha} \boldsymbol{\sigma}^{\alpha} : \frac{\partial \boldsymbol{\epsilon}^{\alpha}}{\partial \boldsymbol{\epsilon}} h^{\alpha}$, which, from $\boldsymbol{\epsilon}^{\alpha} = \boldsymbol{\epsilon} - (\mathbf{a}^{\alpha} \otimes \mathbf{n})_S$, leads to

$$\frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}} = \sum_{\alpha} \boldsymbol{\sigma}^{\alpha} h^{\alpha} - \sum_{\alpha} (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}) \cdot \frac{\partial \mathbf{a}^{\alpha}}{\partial \boldsymbol{\epsilon}} h^{\alpha} = \boldsymbol{\sigma} - (\boldsymbol{\sigma} \cdot \mathbf{n}) \frac{\partial \sum_{\alpha} \mathbf{a}^{\alpha} h^{\alpha}}{\partial \boldsymbol{\epsilon}} = \boldsymbol{\sigma}.$$

Its only drawback - and this is a potentially major one from a physical point of view - is that this is primarily due to the usage of a single common normal vector \mathbf{n} . While the appropriate definition of a normal vector $\mathbf{n}^{\alpha\beta}$ between any phase-pairing is already a topic of some debate, i.e. for example in terms of the $\mathbf{q}^{\alpha\beta}$ or the normalized differences $\nabla \phi^{\beta} - \nabla \phi^{\alpha}$, doing so for all phases and thus also phase-pairings at once is an even more difficult problem. This is also the primary criticism of the model in [62] by Tschukin in [74].

In addition, even though the contribution to the driving force through the dependence on \mathbf{n} is left out in [62], defining \mathbf{n} in terms of a fairly complex averaging procedure based on the phasefield and various mechanical properties will, despite of the simplicity of the expression in Equation (7.172), lead to a very complex total driving force in combination with the corresponding expressions for $\frac{\partial \mathbf{n}}{\partial \phi^\alpha}$ and $\frac{\partial \mathbf{n}}{\partial \nabla \phi^\alpha}$. \diamond

The Model by Schneider et al. ([63] and [61])

A quite different generalization of the model from Section 7.2.2 to the multiphase setting was proposed by Schneider et al. in [63] in the finite deformation setting (see also [61]).

Dropping the restrictions imposed by Equation (7.156) except with respect to a given reference phase R, Equation (7.159) applied for this reference phase becomes

$$\mathbf{F}^R = \mathbf{F} - \sum_{\beta \neq R} \mathbf{a}^{R\beta} \otimes \mathbf{N}^{R\beta} h^{\beta}(\phi) \quad \text{and} \quad \mathbf{F}^{\alpha} = \mathbf{F}^R + \mathbf{a}^{R\alpha} \otimes \mathbf{N}^{R\alpha}, \quad \alpha \neq R \quad (7.173)$$

resp.

$$\boldsymbol{\epsilon}^R = \boldsymbol{\epsilon} - \sum_{\beta \neq R} (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^{\beta}(\phi) \quad \text{and} \quad \boldsymbol{\epsilon}^{\alpha} = \boldsymbol{\epsilon}^R + (\mathbf{a}^{R\alpha} \otimes \mathbf{n}^{R\alpha})_S, \quad \alpha \neq R \quad (7.174)$$

in the small deformation case. The remaining $3(N-1)$ unknowns in the jump vectors $\mathbf{a}^{R\alpha}$, $\alpha \neq R$, can then again be fixed by imposing the continuity conditions on the normal stresses in Equation (7.157) with respect to this reference phase only, i.e.

$$\llbracket \boldsymbol{\sigma} \rrbracket^{R\alpha} \cdot \mathbf{n}^{R\alpha} = (\boldsymbol{\sigma}^\alpha - \boldsymbol{\sigma}^R) \cdot \mathbf{n}^{R\alpha} = \mathbf{0} \quad (7.175)$$

for all $\alpha \neq R$.

Inserting the definition of the $\boldsymbol{\epsilon}^\alpha$ in terms of $\boldsymbol{\epsilon}$ and the $\mathbf{a}^{R\alpha}$ into the phase-specific stress strain relationships, $\alpha = 1, \dots, N$, one has

$$\boldsymbol{\sigma}^\alpha = \begin{cases} \mathbf{C}^R : \left(\boldsymbol{\epsilon} - \sum_{\beta \neq R} (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi) - \tilde{\boldsymbol{\epsilon}}^R \right) & \alpha = R, \\ \mathbf{C}^\alpha : \left(\boldsymbol{\epsilon} - \sum_{\beta \neq R} (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi) + (\mathbf{a}^{R\alpha} \otimes \mathbf{n}^{R\alpha})_S - \tilde{\boldsymbol{\epsilon}}^\alpha \right) & \text{else.} \end{cases}$$

Together with the right subsymmetry of \mathbf{C}^α and defining as before $\boldsymbol{\Sigma}^\alpha := \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$, this simplifies to

$$\boldsymbol{\sigma}^\alpha = \begin{cases} \boldsymbol{\Sigma}^R - \sum_{\beta \neq R} h^\beta(\phi) (\mathbf{C}^R \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} & \alpha = R, \\ \boldsymbol{\Sigma}^\alpha - \sum_{\beta \neq R} h^\beta(\phi) (\mathbf{C}^\alpha \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} + (\mathbf{C}^\alpha \cdot \mathbf{n}^{R\alpha}) \cdot \mathbf{a}^{R\alpha} & \text{else,} \end{cases}$$

and thus, in combination with the condition on the normal continuity of the stresses leads to

$$\left(\left(\boldsymbol{\Sigma}^R - \sum_{\beta \neq R} (\mathbf{C}^R \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} \right) - \left(\boldsymbol{\Sigma}^\alpha - \left(\sum_{\beta \neq R} (\mathbf{C}^\alpha \cdot \mathbf{n}^{R\beta}) \otimes \mathbf{a}^{R\beta} h^\beta \right) + (\mathbf{C}^\alpha \cdot \mathbf{n}^{R\alpha}) \cdot \mathbf{a}^{R\alpha} \right) \right) \cdot \mathbf{n}^{R\alpha} \stackrel{!}{=} \mathbf{0}.$$

Regrouping terms, one therefore has to solve the $3(N-1)$ -system

$$\begin{aligned} & (\mathbf{n}^{R\alpha} \cdot \mathbf{C}^\alpha \cdot \mathbf{n}^{R\alpha}) \cdot \mathbf{a}^{R\alpha} - \sum_{\beta \neq R} h^\beta (\mathbf{n}^{R\alpha} \cdot (\mathbf{C}^\alpha - \mathbf{C}^R) \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} \\ &= \left((\mathbf{n}^{R\alpha} \cdot \mathbf{C}^\alpha \cdot \mathbf{n}^{R\alpha}) - h^\alpha (\mathbf{n}^{R\alpha} \cdot (\mathbf{C}^\alpha - \mathbf{C}^R) \cdot \mathbf{n}^{R\alpha}) \right) \cdot \mathbf{a}^{R\alpha} \\ & \quad - \sum_{\beta \neq R, \alpha} h^\beta (\mathbf{n}^{R\alpha} \cdot (\mathbf{C}^\alpha - \mathbf{C}^R) \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} \quad (7.176) \\ &= (\mathbf{n}^{R\alpha} \cdot ((1 - h^\alpha) \mathbf{C}^\alpha + h^\alpha \mathbf{C}^R) \cdot \mathbf{n}^{R\alpha}) \cdot \mathbf{a}^{R\alpha} - \sum_{\beta \neq R, \alpha} h^\beta (\mathbf{n}^{R\alpha} \cdot (\mathbf{C}^\alpha - \mathbf{C}^R) \cdot \mathbf{n}^{R\beta}) \cdot \mathbf{a}^{R\beta} \\ & \stackrel{!}{=} (\boldsymbol{\Sigma}^R - \boldsymbol{\Sigma}^\alpha) \cdot \mathbf{n}^{R\alpha} = -\llbracket \boldsymbol{\Sigma} \rrbracket^{R\alpha} \cdot \mathbf{n}^{R\alpha} =: \mathbf{r}^{R\alpha} \end{aligned}$$

with the “residual” $\mathbf{r}^{R\alpha}$ corresponding to the (negative) stress jump as it would arise in the Voigt-Taylor model.

Finally, having determined the $\mathbf{a}^{R\alpha}$, the total stress $\boldsymbol{\sigma}$ is obtained, based on a weighted average of the phase-specific ones, as

$$\begin{aligned} \boldsymbol{\sigma} &= \sum_{\alpha=1}^N \boldsymbol{\sigma}^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^N \left(\mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha) \right) h^\alpha(\phi) \\ &= \sum_{\alpha=1}^N \left(\mathbf{C}^\alpha : \left(\boldsymbol{\epsilon} - \sum_{\beta \neq R} (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi) - \tilde{\boldsymbol{\epsilon}}^\alpha \right) + \begin{cases} \mathbf{0} & \alpha = 1, \\ (\mathbf{a}^{R\alpha} \otimes \mathbf{n}^{R\alpha})_S & \text{else} \end{cases} \right) h^\alpha(\phi) \\ &= \sum_{\alpha=1}^N \left(\boldsymbol{\Sigma}^\alpha - \mathbf{C}^\alpha : \left(\sum_{\beta \neq 1} (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi) \right) \right) h^\alpha(\phi) + \sum_{\beta \neq R} \mathbf{C}^\beta : (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi). \end{aligned}$$

Using $1 = \sum_{\alpha=1}^N h^\alpha(\phi)$, this can further be combined to give⁹⁵

$$\boldsymbol{\sigma} = \sum_{\alpha=1}^N \left(\boldsymbol{\Sigma}^\alpha - \left(\sum_{\beta \neq R} (\mathbf{C}^\alpha - \mathbf{C}^\beta) : (\mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S h^\beta(\phi) \right) \right) h^\alpha(\phi). \quad (7.177)$$

Remark 155. Similarly to the chemical case considered in Section 7.1, it is not really necessary to solve for all N phases at once. Instead, the only equations which are actually coupled are the ones for which $h^\alpha(\phi) > 0$, whereas the strains and stresses for all remaining phases can, if required, be recovered a posteriori. \diamond

In contrast to the model based on a common normal, a clear advantage of this model is that it allows, more in line with the standard phasefield approach, using different normal vector $\mathbf{n}^{\alpha\beta}$ for each of the phase-pairings. This comes at a heavy price though, as this is made possible here by enforcing the desired jump-conditions with respect to a single reference phase R only, whereas all other phase-pairings will generally not satisfy the jump conditions on either the strains or the stresses⁹⁶.

Remark 156. Due to this, an important question is also how the reference phase can be chosen “appropriately”. A simple and quite intuitive possibility is to use the phase with the highest local volume-fraction ϕ^α . Other choices are of course possible, i.e. one could for example try to take the mechanical properties of the phases into account. Regardless of the actual choice, the necessity of having to choose is on the one hand an unpleasant one, and on the other hand also entails some issues (there being others as well as will be seen below) in the calculation of the driving force as the points where the reference phases change will generally be associated with discontinuities in the mechanical energies due to the sudden change of the subset of jump conditions to be satisfied. \diamond

Furthermore - and this is an in principle very problematic issue from a variational point of view - even though one can conveniently define an elastic free energy density f_{el} as in Equation (7.71) based on the $\boldsymbol{\epsilon}^\alpha$ in Equation (7.174), there is no extremum principle with respect to this density underlying the actual determination of the jump vectors $\mathbf{a}^{R\alpha}$. In fact, a minimization of the free energy density (7.71) as in [51] for the given form of the $\boldsymbol{\epsilon}^\alpha$ would, with

$$\delta \boldsymbol{\epsilon}^\alpha = \begin{cases} -\sum_{\beta \neq R} h^\beta (\delta \mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S & \alpha = R, \\ -\sum_{\beta \neq R} h^\beta (\delta \mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S + (\delta \mathbf{a}^{R\alpha} \otimes \mathbf{n}^{R\alpha})_S & \alpha \neq R \end{cases}$$

result in

$$\begin{aligned} & -h^R \boldsymbol{\sigma}^R : \sum_{\beta \neq R} h^\beta (\delta \mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S + \sum_{\alpha \neq R} \left(-h^\alpha \boldsymbol{\sigma}^\alpha : \sum_{\beta \neq R} h^\beta (\delta \mathbf{a}^{R\beta} \otimes \mathbf{n}^{R\beta})_S + h^\alpha \boldsymbol{\sigma}^\alpha : (\delta \mathbf{a}^{R\alpha} \otimes \mathbf{n}^{R\alpha})_S \right) \\ & = \sum_{\alpha \neq R} h^\alpha \left(\left(\boldsymbol{\sigma}^\alpha - \sum_{\beta \neq R} h^\beta \boldsymbol{\sigma}^\beta - h^R \boldsymbol{\sigma}^R \right) \cdot \mathbf{n}^{R\alpha} \right) \cdot \delta \mathbf{a}^{R\alpha} = \sum_{\alpha \neq R} h^\alpha \left(\left(\boldsymbol{\sigma}^\alpha - \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta \right) \cdot \mathbf{n}^{R\alpha} \right) \cdot \delta \mathbf{a}^{R\alpha} \stackrel{!}{=} \mathbf{0}, \end{aligned}$$

and thus the requirement that each $\mathbf{a}^{R\alpha}$ should be such that the normal stress components w.r.t. $\mathbf{n}^{R\alpha}$ coincides with the ones of the total resulting stress. As this differs from the characterization in Equation (7.175), there is no reason for the dependence of the $\mathbf{a}^{R\alpha}$ with respect to the parameters ϕ , $\boldsymbol{\epsilon}$ and \mathbf{n} to drop out when differentiating f_{el} . In particular, the stress defined by $\boldsymbol{\sigma} = \sum_{\alpha} \boldsymbol{\sigma}^\alpha h^\alpha = \sum_{\alpha} \mathbf{C}^\alpha : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha)$ does **not** generally satisfy the crucial relation $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$ except within bulk- or two-phase regions.

This has, at least at a theoretical level, far-reaching consequences in terms of the postulated gradient flow for ϕ similar to the discussion in [41]. In fact, one now has to choose between

⁹⁵The summation could clearly also be restricted to $\beta \neq R, \alpha$.

⁹⁶Unfortunately, the claim in [61] that the remaining equations are redundant is incorrect.

maintaining the simple structure of the algorithm, i.e. ignoring the interdependencies between the different parameters ϕ and the displacement respectively ϵ , or, ultimately, the minimization of the functional in Equation (7.67).

Remark 157. From a more practical point of view, this global variational inconsistency can likely be considered as a secondary issue as compared to the many difficult questions related to multiphase regions such as the choice of normal vectors and implications associated to the choice of reference phase and, more importantly, the continued neglect of the dependence on \mathbf{n} in [63].

One can in principle derive a driving force consistent with the minimization of \mathcal{F}_ϵ subject to the momentum-balance based on σ as above considered as a side-constraint using e.g. an adjoint equation corresponding to a reduced formulation in terms of the displacement \mathbf{u} as a function of ϕ and the $\mathbf{n}^{\text{R}\alpha}$. This has the drawback that determining the adjoint state gives rise to a second **global** problem similar in nature to the steady-state momentum-balance Equation (7.69) and is therefore potentially relatively expensive. \diamond

For the calculation of the driving force, it is again convenient to separate the derivative into two separate contributions, one based on fixed normal vectors $(\mathbf{n}^{\alpha\beta})_{\alpha\neq\beta}$ and the other one due to the changes in the normal vectors. A partial differentiation with respect to ϕ^α for fixed normals leads to

$$\frac{\partial}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}}f(\phi, \epsilon) = \sum_{\beta} f^\beta \frac{\partial h^\beta}{\partial\phi^\alpha} + \sum_{\beta} \sigma^\beta : \frac{\partial\epsilon^\beta(\phi, \epsilon)}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} h^\beta(\phi).$$

From

$$\epsilon^{\text{R}}(\phi, \epsilon) = \epsilon - \sum_{\beta\neq\text{R}} (\mathbf{a}^{\text{R}\beta}(\phi, \epsilon) \otimes \mathbf{n}^{\text{R}\beta})_S h^\beta(\phi), \quad \epsilon^\alpha(\phi, \epsilon) = \epsilon^{\text{R}}(\phi, \epsilon) + (\mathbf{a}^{\text{R}\alpha}(\phi, \epsilon) \otimes \mathbf{n}^{\text{R}\alpha})_S \quad \alpha \neq \text{R}$$

together with $\sigma = \sum_{\beta=1}^N \sigma^\beta h^\beta(\phi)$, the second sum can be rewritten to to

$$\begin{aligned} \sum_{\beta} h^\beta \sigma^\beta : \frac{\partial\epsilon^\beta(\phi, \epsilon)}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} &= \left(\sum_{\beta} \sigma^\beta h^\beta \right) : \frac{\partial\epsilon^{\text{R}}(\phi, \epsilon)}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} + \sum_{\beta\neq\text{R}} \sigma^\beta : \left(\frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} \otimes \mathbf{n}^{\text{R}\beta} \right)_S h^\beta(\phi) \\ &= -\sigma : \sum_{\beta\neq\text{R}} (\mathbf{a}^{\text{R}\beta} \otimes \mathbf{n}^{\text{R}\beta})_S \frac{\partial h^\beta}{\partial\phi^\alpha} - \sigma : \sum_{\beta\neq\text{R}} \left(\frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} \otimes \mathbf{n}^{\text{R}\beta} \right)_S h^\beta \\ &\quad + \sum_{\beta\neq\text{R}} \sigma^\beta : \left(\frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha}\Big|_{\mathbf{n}^{\alpha\beta}} \otimes \mathbf{n}^{\text{R}\beta} \right)_S h^\beta \\ &= -\sigma : \sum_{\beta\neq\text{R}} (\mathbf{a}^{\text{R}\beta} \otimes \mathbf{n}^{\text{R}\beta})_S \frac{\partial h^\beta}{\partial\phi^\alpha} + \sum_{\beta\neq\text{R}} \left((\sigma^\beta - \sigma) \cdot \mathbf{n}^{\text{R}\beta} \right) \cdot \frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha} h^\beta. \end{aligned}$$

Using $(\mathbf{a}^{\text{R}\beta}(\phi, \epsilon) \otimes \mathbf{n}^{\text{R}\beta})_S = \epsilon^\beta - \epsilon^{\text{R}}$, an alternative expression is given by

$$\sum_{\beta=1}^N h^\beta \sigma^\beta : \frac{\partial\epsilon^\beta(\phi, \epsilon)}{\partial\phi^\alpha} = -\sigma : \sum_{\beta\neq\text{R}} (\epsilon^\beta - \epsilon^{\text{R}}) \frac{\partial h^\beta}{\partial\phi^\alpha} + \sum_{\beta\neq\text{R}} \left((\sigma^\beta - \sigma) \cdot \mathbf{n}^{\text{R}\beta} \right) \cdot \frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha} h^\beta$$

which, with

$$\sigma : \sum_{\beta\neq\text{R}} \epsilon^{\text{R}} \frac{\partial h^\beta}{\partial\phi^\alpha} = \sigma : \epsilon^{\text{R}} \frac{\partial \sum_{\beta\neq\text{R}} h^\beta}{\partial\phi^\alpha} = \epsilon^{\text{R}} \frac{\partial(1 - h^{\text{R}}(\phi))}{\partial\phi^\alpha} = -\sigma : \epsilon^{\text{R}} \frac{\partial h^{\text{R}}}{\partial\phi^\alpha}$$

show that

$$\sum_{\beta} h^\beta \sigma^\beta : \frac{\partial\epsilon^\beta(\phi, \epsilon)}{\partial\phi^\alpha} = -\sigma : \sum_{\beta=1}^N \epsilon^\beta \frac{\partial h^\beta}{\partial\phi^\alpha} + \sum_{\beta\neq\text{R}} \left((\sigma^\beta - \sigma) \cdot \mathbf{n}^{\text{R}\beta} \right) \cdot \frac{\partial\mathbf{a}^{\text{R}\beta}}{\partial\phi^\alpha} h^\beta.$$

This makes it possible to rewrite this contribution to the total driving force in a manner more similar to the one for the model with the common normals in Equation (7.171) as

$$\frac{\partial}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} f(\phi, \boldsymbol{\epsilon}) = \sum_{\beta} \left(f^\beta - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha} + \sum_{\beta \neq \mathbb{R}} h^\beta \left((\boldsymbol{\sigma}^\beta - \boldsymbol{\sigma}) \cdot \mathbf{n}^{\mathbb{R}\beta} \right) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}}. \quad (7.178)$$

In contrast to the former model, there is no reason here for the terms in $\frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}}$, except of course for two-phase regions where both phases with $h^\beta > 0$ share the only common relevant normal and therefore the normal stress. This is considerably less problematic than the issues caused by $\boldsymbol{\sigma} \neq \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$ as the required expressions for these derivatives can be recovered in a **local** fashion after a differentiation of their defining system in Equation (7.176), leading to

$$\begin{aligned} & \left(\mathbf{n}^{\mathbb{R}\alpha} \cdot \left((1 - h^\alpha) \mathbf{c}^\alpha + h^\alpha \mathbf{c}^{\mathbb{R}} \right) \cdot \mathbf{n}^{\mathbb{R}\alpha} \right) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\alpha}}{\partial \phi^\delta} \Big|_{\mathbf{n}^{\alpha\beta}} - \sum_{\beta \neq \mathbb{R}, \alpha} h^\beta \left(\mathbf{n}^{\mathbb{R}\alpha} \cdot (\mathbf{c}^\alpha - \mathbf{c}^{\mathbb{R}}) \cdot \mathbf{n}^{\mathbb{R}\beta} \right) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \phi^\delta} \Big|_{\mathbf{n}^{\alpha\beta}} \\ &= - \left(\mathbf{n}^{\mathbb{R}\alpha} \cdot \left(\left(1 - \frac{\partial h^\alpha}{\partial \phi^\delta} \right) \mathbf{c}^\alpha + \frac{\partial h^\alpha}{\partial \phi^\delta} \mathbf{c}^{\mathbb{R}} \right) \cdot \mathbf{n}^{\mathbb{R}\alpha} \right) \cdot \mathbf{a}^{\mathbb{R}\alpha} + \sum_{\beta \neq \mathbb{R}, \alpha} \frac{\partial h^\beta}{\partial \phi^\delta} \left(\mathbf{n}^{\mathbb{R}\alpha} \cdot (\mathbf{c}^\alpha - \mathbf{c}^{\mathbb{R}}) \cdot \mathbf{n}^{\mathbb{R}\beta} \right) \cdot \mathbf{a}^{\mathbb{R}\beta}. \end{aligned}$$

Remark 158. A pleasant property of this system is that it is based on the same matrix as the original system (7.176) for the determination of the $\mathbf{a}^{\mathbb{R}\alpha}$ and differs only in the right-hand sides and one can thus reuse a previous factorization for solving for the $\frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}}$. Nevertheless, one needs to solve such a system for the derivative with respect to each of the phases.

A very pragmatic alternative is to instead simply neglect these additional contributions. While this has no real mathematical justification, it leads to a simpler driving force, which remains correct within the two-phase regions, while the errors in the multiphase regions are expected to be of a similar order as the ones already induced by the variational inconsistency due to the construction of $\boldsymbol{\sigma}$. \diamond

As the calculation of the phase-specific strains relies on the $\mathbf{n}^{\mathbb{R}\alpha}$ only, there are no additional contributions to $\frac{\partial f_{el}}{\partial \phi^\alpha}$ and $\frac{\partial f_{el}}{\partial \nabla \phi^\alpha}$ due to the other normal vectors. For the $\mathbf{n}^{\mathbb{R}\alpha}$, one has

$$\begin{aligned} \frac{\partial f_{el}}{\partial \mathbf{n}^{\mathbb{R}\alpha}} &= \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \mathbf{n}^{\mathbb{R}\alpha}} = \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta : \frac{\partial \boldsymbol{\epsilon}^{\mathbb{R}}}{\partial \mathbf{n}^{\mathbb{R}\alpha}} + \sum_{\beta \neq \mathbb{R}} h^\beta \boldsymbol{\sigma}^\beta : \frac{\partial}{\partial \mathbf{n}^{\mathbb{R}\alpha}} (\mathbf{a}^{\mathbb{R}\beta} \otimes \mathbf{n}^{\mathbb{R}\beta})_S \\ &= - \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta : \sum_{\delta \neq \mathbb{R}} h^\delta \frac{\partial}{\partial \mathbf{n}^{\mathbb{R}\alpha}} (\mathbf{a}^{\mathbb{R}\delta} \otimes \mathbf{n}^{\mathbb{R}\delta})_S + \sum_{\beta \neq \mathbb{R}} h^\beta \boldsymbol{\sigma}^\beta : \frac{\partial}{\partial \mathbf{n}^{\mathbb{R}\alpha}} (\mathbf{a}^{\mathbb{R}\beta} \otimes \mathbf{n}^{\mathbb{R}\beta})_S. \end{aligned}$$

Whereas the explicit dependence on the normal vectors only entails a single contribution for $\alpha = \delta$ resp. $\alpha \neq \delta$, the $\mathbf{a}^{\mathbb{R}\delta}$ depend on all $\mathbf{n}^{\mathbb{R}\alpha}$ simultaneously, such that one has

$$\begin{aligned} \frac{\partial f_{el}}{\partial \mathbf{n}^{\mathbb{R}\alpha}} &= - \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta \cdot \mathbf{a}^{\mathbb{R}\alpha} h^\alpha - \sum_{\beta} \sum_{\delta \neq \mathbb{R}} h^\delta h^\beta (\boldsymbol{\sigma}^\beta \cdot \mathbf{n}^{\mathbb{R}\delta}) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\delta}}{\partial \mathbf{n}^{\mathbb{R}\alpha}} \\ &\quad + h^\alpha \boldsymbol{\sigma}^\alpha \cdot \mathbf{a}^{\mathbb{R}\alpha} + \sum_{\beta \neq \mathbb{R}} h^\beta (\boldsymbol{\sigma}^\beta \cdot \mathbf{n}^{\mathbb{R}\beta}) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \mathbf{n}^{\mathbb{R}\alpha}}. \end{aligned}$$

Using $\boldsymbol{\sigma} = \sum_{\beta} h^\beta \boldsymbol{\sigma}^\beta$ and exchanging the dummy-indices, this can be rewritten in a more compact form as

$$\frac{\partial f_{el}}{\partial \mathbf{n}^{\mathbb{R}\alpha}} = h^\alpha (\boldsymbol{\sigma}^\alpha - \boldsymbol{\sigma}) \cdot \mathbf{a}^{\mathbb{R}\alpha} + \sum_{\beta \neq \mathbb{R}} h^\beta \left((\boldsymbol{\sigma}^\beta - \boldsymbol{\sigma}) \cdot \mathbf{n}^{\mathbb{R}\beta} \right) \cdot \frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \mathbf{n}^{\mathbb{R}\alpha}}. \quad (7.179)$$

There is again no reason for the second contribution to drop out. Unless one decides to neglect this contribution for simplicity, the required contributions through $\frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \mathbf{n}^{\mathbb{R}\alpha}}$ can be recovered, similarly to those due to $\frac{\partial \mathbf{a}^{\mathbb{R}\beta}}{\partial \phi^\alpha}$ through a linearization of the system (7.176).

The Model by Tschukin ([74])

A quite different extension of the two-phase models from Section 7.2.2 is proposed by Tschukin in [74]. In contrast to the work by Schneider, the model in [74] is a priori not based upon a prefixed expression of the strains and stresses in each phase but instead obtained through a generalization of the expression (7.151) of the elastic free energy density. He suggests replacing the first part in Equation (7.151) through a summation over the analogous expression over all phases and, in analogy to the treatment of the surface energy contributions through a summation over all phase-pairings in Section 6.1, replacing the last term through a summation over all two-phase pairings. More precisely, using the same definitions $\bar{\mathcal{S}}_{nn}^{\alpha\beta} = \bar{\mathcal{S}}_{nn}^{\beta\alpha} = (\bar{\mathcal{C}}_{nn}^{\alpha\beta})^\dagger$ where $\bar{\mathcal{C}}_{nn}^{\alpha\beta} = \mathcal{N} : (h^\alpha \mathcal{C}^\beta + h^\beta \mathcal{C}^\alpha) : \mathcal{N}$, he introduces the expression⁹⁷

$$f_{el}(\phi, \nabla\phi, \epsilon) = \frac{1}{2} \sum_{\alpha=1}^N \left(\Sigma^\alpha : (\epsilon - \tilde{\epsilon}^\alpha) - \frac{1}{2} \sum_{\beta \neq \alpha} \llbracket \Sigma \rrbracket^{\alpha\beta} : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : \llbracket \Sigma \rrbracket^{\alpha\beta} h^\beta(\phi) \right) h^\alpha(\phi) \quad (7.180)$$

for the elastic free energy density. The average (effective) stress σ is then derived from Equation (7.180) in accordance to Equation (7.70) as $\sigma := \frac{\partial f_{el}}{\partial \epsilon}$, which, together with

$$\frac{\partial}{\partial \epsilon} \llbracket \Sigma \rrbracket^{\alpha\beta} = \frac{\partial}{\partial \epsilon} (\mathcal{C}^\beta : (\epsilon - \tilde{\epsilon}^\beta) - \mathcal{C}^\alpha : (\epsilon - \tilde{\epsilon}^\alpha)) = \mathcal{C}^\beta - \mathcal{C}^\alpha$$

leads to the expression

$$\sigma := \frac{\partial f_{el}}{\partial \epsilon} = \sum_{\alpha=1}^N \left(\Sigma^\alpha - \frac{1}{2} \sum_{\beta \neq \alpha} (\mathcal{C}^\beta - \mathcal{C}^\alpha) : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : \llbracket \Sigma \rrbracket^{\alpha\beta} h^\beta(\phi) \right) h^\alpha(\phi) \quad (7.181)$$

generalizing the one in Equation (7.149) from the two-phase case.

Remark 159. A clear advantages as compared to the model by [63], [61] is that it is obvious from equations (7.180) and (7.181) that there is no “preferential” or reference phase. In addition, the stress tensor will, by its very definition in Equation (7.181), satisfy the consistency condition (7.70). \diamond

In contrast to the model by Schneider, Tschukin avoids the issue of defining phase-specific strains and stresses. In particular, the free energy density being based upon considering a sum over essentially “independent” two-phase interactions, it is not at all clear a priori how his model can be extended to situations (such as viscoelastic or plastic problems), where one would require a consistent (i.e. independent of the phase-pairing) definition of these phase-inherent quantities. Simply **defining** similar to the discussion in the two-phase case in Section 7.2.3 the jump vector $\mathbf{a}^{\alpha\beta}$ for a given phase-pairing through $(\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S = \bar{\mathcal{S}}_{nn}^{\alpha\beta} : (\Sigma^\alpha - \Sigma^\beta) = -\bar{\mathcal{S}}_{nn}^{\alpha\beta} : \llbracket \Sigma \rrbracket^{\alpha\beta}$, or, in a simpler form similar to Lemma 9, through

$$\mathbf{a}^{\alpha\beta} = (\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta})^{-1} \cdot ((\Sigma^\alpha - \Sigma^\beta) \cdot \mathbf{n}^{\alpha\beta}) = -(\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta})^{-1} \cdot (\llbracket \Sigma \rrbracket^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}), \quad (7.182)$$

with $\bar{\mathcal{C}}^{\alpha\beta} = h^\beta \mathcal{C}^\alpha + h^\alpha \mathcal{C}^\beta$, it turns out that it is in fact possible to obtain the stress-strain relationship above based upon consistently defined phase-specific stresses σ^α and strains ϵ^α given by

$$\epsilon^\alpha = \epsilon - \sum_{\beta \neq \alpha} h^\beta (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S \quad \text{and} \quad \sigma^\alpha = \mathcal{C}^\alpha : (\epsilon^\alpha - \tilde{\epsilon}^\alpha). \quad (7.183)$$

⁹⁷Note that the additional factor $\frac{1}{2}$ in front of the last term is due to each phase-pairing appearing twice in this generalization and can be eliminated by including each pairing only once by e.g. instead summing over all $1 \leq \alpha < \beta \leq N$.

In fact, using $\bar{\mathcal{S}}_{nn}^{\alpha\beta} = \bar{\mathcal{S}}_{nn}^{\beta\alpha}$ and $[[\Sigma]]^{\alpha\beta} = -[[\Sigma]]^{\beta\alpha}$, the second term in Equation (7.181) can be rewritten as

$$\begin{aligned} & \frac{1}{2} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \mathcal{C}^\alpha : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\beta(\phi) h^\alpha(\phi) - \frac{1}{2} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \mathcal{C}^\beta : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\beta(\phi) h^\alpha(\phi) \\ &= \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \mathcal{C}^\alpha : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\beta(\phi) h^\alpha(\phi), \end{aligned}$$

thus leading to the alternative expression

$$\boldsymbol{\sigma} := \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}} = \sum_{\alpha=1}^N \left(\boldsymbol{\Sigma}^\alpha + \mathcal{C}^\alpha : \left(\sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\beta(\phi) \right) \right) h^\alpha(\phi). \quad (7.184)$$

for the effective stress. As $\boldsymbol{\Sigma}^\alpha = \mathcal{C}^\alpha : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha)$, this is the same as

$$\begin{aligned} \boldsymbol{\sigma} &= \sum_{\alpha=1}^N \mathcal{C}^\alpha : \left(\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha + \left(\sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\beta(\phi) \right) \right) h^\alpha(\phi) \\ &= \sum_{\alpha=1}^N \mathcal{C}^\alpha : \left(\boldsymbol{\epsilon} - \sum_{\beta \neq \alpha} h^\beta (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S - \tilde{\boldsymbol{\epsilon}}^\alpha \right) h^\alpha(\phi), \end{aligned}$$

and thus a weighted average of the phase-specific stresses defined as in Equation (7.183). Furthermore, $\boldsymbol{\epsilon}$ and the $\boldsymbol{\epsilon}^\alpha$ also satisfy the averaging condition $\boldsymbol{\epsilon} = \sum_{\alpha} \boldsymbol{\epsilon}^\alpha h^\alpha(\phi)$, since

$$\sum_{\alpha=1}^N \boldsymbol{\epsilon}^\alpha h^\alpha(\phi) = \sum_{\alpha=1}^N \boldsymbol{\epsilon} h^\alpha(\phi) - \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S h^\alpha(\phi) h^\beta(\phi)$$

and the last term vanishes due to

$$\begin{aligned} & - \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S h^\alpha(\phi) h^\beta(\phi) = \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nn}^{\alpha\beta} : (\boldsymbol{\Sigma}^\beta - \boldsymbol{\Sigma}^\alpha) h^\alpha(\phi) h^\beta(\phi) \\ &= \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nn}^{\alpha\beta} : \boldsymbol{\Sigma}^\beta h^\alpha(\phi) h^\beta(\phi) - \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nn}^{\alpha\beta} : \boldsymbol{\Sigma}^\alpha h^\alpha(\phi) h^\beta(\phi) \end{aligned}$$

and thus

$$\sum_{\alpha=1}^N \sum_{\beta \neq \alpha} (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S h^\alpha(\phi) h^\beta(\phi) = \mathbf{0}.$$

The model in [74] therefore avoids two of the major drawbacks of the model in [63], namely the dependence on a reference phase and the inconsistency between $\boldsymbol{\sigma}$ and $\frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$. Unfortunately, it has other issues within multiphase regions. Firstly, since the calculation of the jump-vectors $\mathbf{a}^{\alpha\beta}$ is such that it would only enforce the normal continuity of the stresses if α and β were the only phases present, there will generally be no phase-pairing for which the phase-specific stresses defined in Equation (7.183) will actually satisfy this continuity condition in multiphase regions. Secondly, and this is probably the more problematic point, the definition of f_{el} in Equation (7.180) is much less natural than it may seem at first sight. More precisely, with

$$\begin{aligned} & - \frac{1}{4} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} (\boldsymbol{\Sigma}^\beta - \boldsymbol{\Sigma}^\alpha) : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\alpha(\phi) h^\beta(\phi) \\ &= \frac{1}{4} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \boldsymbol{\Sigma}^\alpha : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\alpha(\phi) h^\beta(\phi) + \frac{1}{4} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \boldsymbol{\Sigma}^\beta : \bar{\mathcal{S}}_{nn}^{\beta\alpha} : [[\Sigma]]^{\beta\alpha} h^\alpha(\phi) h^\beta(\phi) \\ &= \frac{1}{2} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \boldsymbol{\Sigma}^\alpha : \bar{\mathcal{S}}_{nn}^{\alpha\beta} : [[\Sigma]]^{\alpha\beta} h^\alpha(\phi) h^\beta(\phi) \end{aligned}$$

f_{el} can be written as

$$f_{el}(\phi, \nabla\phi, \epsilon) = \frac{1}{2} \sum_{\alpha=1}^N \left(\Sigma^\alpha : (\epsilon - \tilde{\epsilon}^\alpha) + \Sigma^\alpha : \left(\sum_{\beta \neq \alpha} \bar{\mathcal{S}}_{nm}^{\alpha\beta} : \llbracket \Sigma \rrbracket^{\alpha\beta} h^\beta(\phi) \right) \right) h^\alpha(\phi).$$

Together the definition of the $\mathbf{a}^{\alpha\beta}$ and the ϵ^α in Equation (7.183), this further simplifies to

$$f_{el}(\phi, \nabla\phi, \epsilon) = \frac{1}{2} \sum_{\alpha=1}^N \Sigma^\alpha(\epsilon) : (\epsilon^\alpha(\phi, \nabla\phi, \epsilon) - \tilde{\epsilon}^\alpha) h^\alpha(\phi), \quad (7.185)$$

and thus a free energy density where one half of the expression is based on the supposedly bad approximation $\Sigma^\alpha = \mathbf{C}^\alpha : (\epsilon - \tilde{\epsilon}^\alpha)$ instead of σ^α corresponding to the ‘‘correct’’ strains ϵ^α . One particularly problematic point with this definition is that, whereas the weighted average $\frac{1}{2} \sum_{\alpha} (\epsilon^\alpha - \tilde{\epsilon}^\alpha) : \mathbf{C}^\alpha : (\epsilon^\alpha - \tilde{\epsilon}^\alpha) h^\alpha(\phi)$ is necessarily non-negative regardless of the precise choice of the ϵ^α , this need not be the case for the expression in Equation (7.185), and it can in fact be observed numerically that defining f_{el} as in [74] may indeed lead to locally negative free energy densities.

The last expression in Equation (7.185) also allows for a somewhat simpler derivation of the correct driving force contributions (in [74] based on the $\mathbf{n}^{\alpha\beta}$ defined in terms of the gradients of ϕ only)

$$\frac{\partial f_{el}}{\partial \phi^\alpha} = \frac{1}{2} \sum_{\beta=1}^N \left(\Sigma^\beta : (\epsilon - \tilde{\epsilon}^\beta) - \sum_{\delta \neq \beta} (\mathbf{a}^{\beta\delta} \otimes \mathbf{n}^{\beta\delta})_S : \mathbf{C}^\beta : (\mathbf{a}^{\beta\delta} \otimes \mathbf{n}^{\beta\delta})_S (h^\delta)^2 \right) \frac{\partial h^\beta}{\partial \phi^\alpha} \quad (7.186)$$

as in [74], which, while correct, requires a somewhat heavy formalism due to the use of the pseudo-inverses in Equation (7.180). A partial differentiation with respect to ϕ^α while holding the normal vectors fixed first leads to

$$\left. \frac{\partial f_{el}}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}} = \frac{1}{2} \sum_{\beta} \Sigma^\beta : (\epsilon^\beta - \tilde{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha} + \frac{1}{2} \sum_{\beta} h^\beta \Sigma^\beta : \left. \frac{\partial \epsilon^\beta}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}},$$

where, by Equation (7.183),

$$\left. \frac{\partial \epsilon^\beta}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}} = - \sum_{\delta \neq \beta} (\mathbf{a}^{\beta\delta} \otimes \mathbf{n}^{\beta\delta})_S \frac{\partial h^\delta}{\partial \phi^\alpha} - \sum_{\delta \neq \beta} h^\delta \left(\left. \frac{\partial \mathbf{a}^{\beta\delta}}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}} \otimes \mathbf{n}^{\beta\delta} \right)_S.$$

By a partial differentiation of the defining Equation (7.182) for the jump vectors and the definition of $\mathbf{C}^{\beta\delta}$, it is easily seen that

$$\left. \frac{\partial \mathbf{a}^{\beta\delta}}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}} = - (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot \left(\frac{\partial h^\delta}{\partial \phi^\alpha} \mathbf{C}^\beta + \frac{\partial h^\beta}{\partial \phi^\alpha} \mathbf{C}^\delta \right) \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta}.$$

Further expanding the term $\Sigma^\beta : (\epsilon^\beta - \tilde{\epsilon}^\beta)$ using $\epsilon^\beta = \epsilon - \sum_{\delta \neq \beta} (\mathbf{a}^{\beta\delta} \otimes \mathbf{n}^{\beta\delta})_S h^\delta$, it follows that

$$\begin{aligned} \left. \frac{\partial f_{el}}{\partial \phi^\alpha} \right|_{\mathbf{n}^{\alpha\beta}} &= \frac{1}{2} \sum_{\beta} \Sigma^\beta : (\epsilon - \tilde{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha} \\ &\quad - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} (\Sigma^\beta \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} h^\delta \frac{\partial h^\beta}{\partial \phi^\alpha} - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^\beta (\Sigma^\beta \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^\delta}{\partial \phi^\alpha} \\ &\quad + \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^\beta h^\delta (\Sigma^\beta \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot \left(\frac{\partial h^\delta}{\partial \phi^\alpha} \mathbf{C}^\beta + \frac{\partial h^\beta}{\partial \phi^\alpha} \mathbf{C}^\delta \right) \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \end{aligned}$$

Rewriting the last line as

$$\begin{aligned}
& \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\delta} \mathbf{C}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\delta}}{\partial \phi^{\alpha}} \\
& + \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} h^{\delta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\beta} \mathbf{C}^{\delta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} \\
& = \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot (h^{\delta} \mathbf{C}^{\beta} + h^{\beta} \mathbf{C}^{\delta} - h^{\beta} \mathbf{C}^{\delta}) \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\delta}}{\partial \phi^{\alpha}} \\
& + \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} h^{\delta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot (h^{\delta} \mathbf{C}^{\beta} + h^{\beta} \mathbf{C}^{\delta} - h^{\delta} \mathbf{C}^{\beta}) \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} \\
& = \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\delta}}{\partial \phi^{\alpha}} + \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\delta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} \\
& - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\beta} \mathbf{C}^{\delta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\delta}}{\partial \phi^{\alpha}} \\
& - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\delta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\delta} \mathbf{C}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}
\end{aligned}$$

allows canceling the middle row, leaving

$$\begin{aligned}
\left. \frac{\partial f_{el}}{\partial \phi^{\alpha}} \right|_{\mathbf{n}^{\alpha\beta}} &= \frac{1}{2} \sum_{\beta} \boldsymbol{\Sigma}^{\beta} : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\beta}) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} \\
& - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\beta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\beta} \mathbf{C}^{\delta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\delta}}{\partial \phi^{\alpha}} \\
& - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\delta} (\boldsymbol{\Sigma}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\delta} \mathbf{C}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}.
\end{aligned}$$

Finally, exchanging dummy-indices and making use of $\mathbf{n}^{\beta\delta} = -\mathbf{n}^{\delta\beta}$ and $\mathbf{a}^{\beta\delta} = \mathbf{a}^{\delta\beta}$, the last two rows can further be summarized to

$$-\frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} h^{\delta} \left((\boldsymbol{\Sigma}^{\beta} - \boldsymbol{\Sigma}^{\delta}) \cdot \mathbf{n}^{\beta\delta} \right) \cdot (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot (\mathbf{n}^{\beta\delta} \cdot h^{\delta} \mathbf{C}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}},$$

which, since $\mathbf{a}^{\beta\delta} = (\mathbf{n}^{\beta\delta} \cdot \bar{\mathbf{C}}^{\beta\delta} \cdot \mathbf{n}^{\beta\delta})^{-1} \cdot \left((\boldsymbol{\Sigma}^{\beta} - \boldsymbol{\Sigma}^{\delta}) \cdot \mathbf{n}^{\beta\delta} \right)$, leads to the relatively simple final expression

$$\left. \frac{\partial f_{el}}{\partial \phi^{\alpha}} \right|_{\mathbf{n}^{\alpha\beta}} = \frac{1}{2} \sum_{\beta} \boldsymbol{\Sigma}^{\beta} : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\beta}) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} - \frac{1}{2} \sum_{\beta} \sum_{\delta \neq \beta} (h^{\delta})^2 \mathbf{a}^{\beta\delta} \cdot (\mathbf{n}^{\beta\delta} \cdot \mathbf{C}^{\beta} \cdot \mathbf{n}^{\beta\delta}) \cdot \mathbf{a}^{\beta\delta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}$$

which amounts to the same as the one in Equation (7.186) by the subsymmetries of \mathbf{C}^{β} .

For the remaining contributions arising from the dependence on the $\mathbf{n}^{\alpha\beta}$, it is convenient to explicitly make use of $\mathbf{n}^{\alpha\beta} = -\mathbf{n}^{\beta\alpha}$ in order to reduce f_{el} to a function of e.g. the $(\mathbf{n}^{\alpha\beta})_{\alpha < \beta}$ only. Doing so and with $\mathbf{a}^{\alpha\beta} = \mathbf{a}^{\beta\alpha}$, it follows that

$$\begin{aligned}
\frac{\partial f_{el}}{\partial \mathbf{n}^{\alpha\beta}} &= \frac{1}{2} \sum_{\delta} h^{\delta} \boldsymbol{\Sigma}^{\delta} : \frac{\partial \boldsymbol{\epsilon}^{\delta}}{\partial \mathbf{n}^{\alpha\beta}} = \frac{1}{2} \sum_{\delta} h^{\delta} \boldsymbol{\Sigma}^{\delta} : \frac{\partial}{\partial \mathbf{n}^{\alpha\beta}} \left(\boldsymbol{\epsilon} - \sum_{\eta \neq \delta} h^{\eta} (\mathbf{a}^{\delta\eta} \otimes \mathbf{n}^{\delta\eta})_S \right) \\
&= -\frac{1}{2} h^{\alpha} h^{\beta} (\boldsymbol{\Sigma}^{\alpha} - \boldsymbol{\Sigma}^{\beta}) \cdot \mathbf{a}^{\alpha\beta} - \frac{1}{2} h^{\alpha} h^{\beta} \left((\boldsymbol{\Sigma}^{\alpha} - \boldsymbol{\Sigma}^{\beta}) \cdot \mathbf{n}^{\alpha\beta} \right) \cdot \frac{\partial \mathbf{a}^{\alpha\beta}}{\partial \mathbf{n}^{\alpha\beta}}
\end{aligned}$$

Since $d\mathbf{a}^{\alpha\beta}$ is, for fixed ϕ , characterized in terms of $\delta\mathbf{n}^{\alpha\beta}$ as

$$(\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}) \cdot d\mathbf{a}^{\alpha\beta} = (\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot d\mathbf{n}^{\alpha\beta} - (d\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}) \cdot \mathbf{a}^{\alpha\beta} - (\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot d\mathbf{n}^{\alpha\beta}) \cdot \mathbf{a}^{\alpha\beta},$$

one has

$$\begin{aligned} ((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot \mathbf{n}^{\alpha\beta}) \cdot d\mathbf{a}^{\alpha\beta} &= ((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot \mathbf{n}^{\alpha\beta}) \cdot (\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta})^{-1} \\ &\quad \cdot \left((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot d\mathbf{n}^{\alpha\beta} - (d\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{n}^{\alpha\beta}) \cdot \mathbf{a}^{\alpha\beta} - (\mathbf{n}^{\alpha\beta} \cdot \bar{\mathcal{C}}^{\alpha\beta} \cdot d\mathbf{n}^{\alpha\beta}) \cdot \mathbf{a}^{\alpha\beta} \right) \\ &= \mathbf{a}^{\alpha\beta} \cdot \left((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot d\mathbf{n}^{\alpha\beta} - 2 \left((\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta} \right)_S : \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{a}^{\alpha\beta} \right) \cdot d\mathbf{n}^{\alpha\beta}, \end{aligned}$$

where use was made of the symmetries of $\bar{\mathcal{C}}^{\alpha\beta}$. Combining this with the previous expression for $\frac{\partial f_{el}}{\partial \mathbf{n}^{\alpha\beta}}$, one therefore has

$$\frac{\partial f_{el}}{\partial \mathbf{n}^{\alpha\beta}} = -h^\alpha h^\beta \left((\boldsymbol{\Sigma}^\alpha - \boldsymbol{\Sigma}^\beta) \cdot \mathbf{a}^{\alpha\beta} - (\mathbf{a}^{\alpha\beta} \otimes \mathbf{n}^{\alpha\beta})_S : \bar{\mathcal{C}}^{\alpha\beta} \cdot \mathbf{a}^{\alpha\beta} \right), \quad (7.187)$$

which is a slightly simpler version of expressions in [74] since it avoids the use of the derivative of the projector \mathcal{N} onto the normal subspace of symmetric second-order tensors.

Remark 160. Summarizing the discussion of the three models considered here, each has its own advantages and disadvantages. Whereas the model in [62] has the mathematically most pleasant properties and in particular satisfies the desired jump conditions between all phase-pairings, this can only be achieved through a very significant geometrical simplification of using a common normal vector for all phase-pairings. The models in [63] and [74] do not make this assumption, but are therefore not generally able to satisfy all jump conditions simultaneously except for in the two-phase regions.

Whereas the model in [63] has the advantage of enforcing the jump conditions on both the strains and stresses at least with respect to the reference phase R, the use of such a preferential phase is clearly an undesirable feature. In contrast, the model in [74] does not rely on a reference phase, but will generally not satisfy any jump conditions on either the stresses or strains within multiphase regions.

From an energetic point of view, both models have some deficiencies when more than two phases are present at a given point. Whereas the small-deformation analogon of the model in [63] is seemingly based on the very natural definition f_{el} as a weighted average of the elastic free energy densities $f_{el}^\alpha(\boldsymbol{\epsilon}^\alpha)$ evaluated in the respective phase-specific strain and thus in particular satisfies $\boldsymbol{\sigma}^\alpha = \frac{\partial f_{el}^\alpha}{\partial \boldsymbol{\epsilon}^\alpha}$, it does generally not satisfy $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$. In contrast, the effective stress in [74] satisfies this relation by construction, but in the form in Equation (7.180) does not allow for such an easy interpretation in terms of phase-specific free energy densities. Even though it was seen in Equation (7.185) that one can artificially rewrite f_{el} in such a form, the resulting expression is quite unsatisfactory as it is based on a mixture of a Voigt-Taylor-type stress prediction and a more quantitative prediction of the strains and does in particular not satisfy $\boldsymbol{\sigma}^\alpha = \frac{\partial f_{el}^\alpha}{\partial \boldsymbol{\epsilon}^\alpha}$.

This should not be misinterpreted as saying that these are bad models, as both achieve a fairly reasonable extension of the very satisfactory two-phase model to within the highly difficult multiphase. Nevertheless, their drawbacks have to be kept in mind when choosing between them for a particular application, and will be inherited from any extension to more complex physical situations. \diamond

7.2.5 Chemo-Elasticity

A more complex situation arises when there is a coupling between the both chemical and elastic influences and the evolution of a microstructure. Such a situation can arise in a diffusion process under a simultaneous elastic deformations if the elastic free energy density depends not only upon the total strain and a phasedependent eigenstrain, but in addition on the concentration itself, for example if the eigenstrains and/or the stiffness tensors have an additional dependence upon the concentration vector. In the bulk phases, this leads to a total free energy densities depending both on the “purely” chemical part $f_{ch}^\alpha(\mathbf{c}^\alpha)$ as well as an additional elastic contribution $f_{el}^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha)$, i.e.

$$f^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha) = f_{ch}^\alpha(\mathbf{c}^\alpha) + f_{el}^\alpha(\boldsymbol{\epsilon}^\alpha, \mathbf{c}^\alpha) = f_{ch}^\alpha(\mathbf{c}^\alpha) + \frac{1}{2}(\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c}^\alpha)) : \mathbf{C}^\alpha(\mathbf{c}^\alpha) : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c}^\alpha)). \quad (7.188)$$

For constructing a phasefield model based on these free energy densities, one again has to choose how this definition should be extended to within the diffuse interface region, where several phases coexist. Assuming as before that one wishes to formulate the model purely in terms of the average concentration \mathbf{c} and strain $\boldsymbol{\epsilon}$, assumed to satisfy

$$\mathbf{c} = \sum_{\alpha=1}^N \mathbf{c}^\alpha h^\alpha(\phi) \quad \text{and} \quad \boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^\alpha h^\alpha(\phi), \quad (7.189)$$

an interpolation of the phase-specific energies $f^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha)$ in Equation (7.188) of the form

$$f(\phi, \mathbf{c}, \boldsymbol{\epsilon}) = \sum_{\alpha=1}^N f^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha) h^\alpha(\phi) \quad (7.190)$$

requires an appropriate additional set of conditions for the determination of the phase-specific quantities $(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha)_{1 \leq \alpha \leq N}$ in terms of ϕ, \mathbf{c} and $\boldsymbol{\epsilon}$.

Probably the simplest approach is to combine the assumption $\mathbf{c}^\alpha = \mathbf{c}, \alpha = 1, \dots, N$, with one of the classical mechanical models, i.e. the Voigt-Taylor or Reuss-Sachs approach. In the former case, the assumption of the equality of all strains $\boldsymbol{\epsilon}^\alpha = \boldsymbol{\epsilon}, \alpha = 1, \dots, N$ leads to the total free energy density

$$f(\phi, \mathbf{c}, \boldsymbol{\epsilon}) = f_{ch}(\phi, \mathbf{c}) + \sum_{\alpha=1}^N \left(\frac{1}{2}(\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) : \mathbf{C}^\alpha(\mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) \right) h^\alpha(\phi).$$

The major advantage of this model lies in its particularly simple form, as, assuming the individual parts are explicit in \mathbf{c} , it results in the formula

$$\boldsymbol{\sigma} = \frac{\partial f}{\partial \boldsymbol{\epsilon}}(\phi, \mathbf{c}, \boldsymbol{\epsilon}) = \sum_{\alpha=1}^N \left(\mathbf{C}^\alpha(\mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) \right) h^\alpha(\phi)$$

for the average stress and

$$\boldsymbol{\mu} := \frac{\partial f}{\partial \mathbf{c}}(\phi, \mathbf{c}, \boldsymbol{\epsilon}) = \underbrace{\frac{\partial f_{ch}}{\partial \mathbf{c}}(\phi, \mathbf{c})}_{=:\boldsymbol{\mu}_{ch}(\mathbf{c})} + \underbrace{\sum_{\alpha=1}^N \left(- \left(\mathbf{C}^\alpha(\mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) \right) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}} + \frac{1}{2}(\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) : \left((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c})) : \frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}} \right) \right)}_{=:\boldsymbol{\mu}_{el}(\mathbf{c}, \boldsymbol{\epsilon})} h^\alpha(\phi)$$

for the chemical potential, both of which are “fully explicit” in \mathbf{c} and $\boldsymbol{\epsilon}$. In the latter case, the situation is - similarly to before - slightly more complex due to the appearance of the derivative of an inverse matrix. As for Equation (7.82), the assumption of equal stresses $\boldsymbol{\sigma}^\alpha = \boldsymbol{\sigma} \forall \alpha$ leads to the total free energy density

$$f(\phi, \mathbf{c}, \boldsymbol{\epsilon}) = f_{ch}(\phi, \mathbf{c}) + \frac{1}{2} \left((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\phi, \mathbf{c})) : \mathbf{C}_{RS}(\phi, \mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\phi, \mathbf{c})) \right), \quad \tilde{\boldsymbol{\epsilon}}(\phi, \mathbf{c}) = \sum_{\alpha=1}^N \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c}) h^\alpha(\phi)$$

with $\mathbf{C}_{RS}(\boldsymbol{\phi}, \mathbf{c}) = \left(\sum_{\alpha=1}^N \mathbf{S}^{\alpha}(\mathbf{c}) h^{\alpha}(\boldsymbol{\phi}) \right)^{-1}$. From this, the stress calculation is in fact even slightly simpler than for the Voigt-Talyer model as it now only involves the averaged eigenstrain $\tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, \mathbf{c}) = \sum_{\alpha=1}^N \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}) h^{\alpha}(\boldsymbol{\phi})$,

$$\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) = \mathbf{C}_{RS}(\boldsymbol{\phi}, \mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, \mathbf{c})).$$

The definition of $\mathbf{C}_{RS}(\boldsymbol{\phi}, \mathbf{c})$ in terms of the inverse of the interpolated compliances gives rise to the same minor complication already encountered in the calculation of the drivingforce for the Reuss-Sachs model as

$$\boldsymbol{\mu}_{el}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) = \frac{\partial f_{el}}{\partial \mathbf{c}}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) = -\mathbf{C}_{RS}(\boldsymbol{\phi}, \mathbf{c}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\mathbf{c})) : \frac{\partial \tilde{\boldsymbol{\epsilon}}}{\partial \mathbf{c}}(\boldsymbol{\phi}, \mathbf{c}) + \frac{1}{2}(\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, \mathbf{c})) : \left((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\mathbf{c})) : \frac{\partial \mathbf{C}_{RS}(\boldsymbol{\phi}, \mathbf{c})}{\partial \mathbf{c}} \right).$$

By the same argument as for Equation (7.84) with \mathbf{c} taking the role of $\boldsymbol{\phi}$, this can again be replaced by a fully explicit expression in terms of the stresses and the derivatives of the $\mathbf{S}^{\alpha}(\mathbf{c})$ leading to

$$\boldsymbol{\mu}_{el}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) = -\boldsymbol{\sigma}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}}{\partial \mathbf{c}}(\boldsymbol{\phi}, \mathbf{c}) + \frac{1}{2} \boldsymbol{\sigma}(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) : \left(\sum_{\alpha=1}^N \frac{\partial \mathbf{S}^{\alpha}(\mathbf{c})}{\partial \mathbf{c}} h^{\alpha}(\boldsymbol{\phi}) \right) \right).$$

It is clear that, while this model is simple, it will inherit the same drawbacks already discussed in the chemical and mechanical case above. As a natural extension of the discussion in Section 7.1, one could instead consider a “grandchem”-type approach with the phase-specific concentrations \mathbf{c}^{α} fixed through equality chemical potentials, i.e.

$$f(\boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon}) = \sum_{\alpha=1}^N \left(f_{ch}^{\alpha}(\mathbf{c}^{\alpha}) + \frac{1}{2} (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \mathbf{C}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) \right) h^{\alpha}(\boldsymbol{\phi})$$

with $(\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}$ fixed by either the equality up to a phase-specific multiplier $\lambda^{\alpha} \mathbf{e}$ of the chemical potentials

$$\boldsymbol{\mu}^{\alpha} = \frac{\partial f_{ch}^{\alpha}}{\partial \mathbf{c}^{\alpha}} - \left(\mathbf{C}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) \right) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^{\alpha}}{\partial \mathbf{c}^{\alpha}} + \frac{1}{2} (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \left((\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \frac{\partial \mathbf{C}^{\alpha}}{\partial \mathbf{c}^{\alpha}} \right) \quad (7.191)$$

or the equality of the $\tilde{\boldsymbol{\mu}}^{\alpha}$ when using a reduced formulation, where the first purely chemical distribution will be designated by $\boldsymbol{\mu}_{ch}^{\alpha}$ and the second one by $\boldsymbol{\mu}_{el}^{\alpha}$. The complexity of such a model depends heavily upon the particular mechanical model this is combined with.

The Voigt-Taylor Case

The simplest situation arises if the $\boldsymbol{\epsilon}^{\alpha}$ are, as in the Voigt-Taylor model, all assumed to be equal to the average strain $\boldsymbol{\epsilon}$. Even though the determination of the mechanical stress does then require the knowledge of the $\tilde{\boldsymbol{\epsilon}}^{\alpha}$, the only unknown parameter for the evaluation of the chemical potentials are the \mathbf{c}^{α} ,

$$\boldsymbol{\mu}^{\alpha} = \boldsymbol{\mu}_{ch}^{\alpha} - \left(\mathbf{C}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) \right) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^{\alpha}}{\partial \mathbf{c}^{\alpha}} + \frac{1}{2} (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \left((\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \frac{\partial \mathbf{C}^{\alpha}}{\partial \mathbf{c}^{\alpha}} \right).$$

One can therefore proceed in a two-step fashion by first solving for the $\tilde{\boldsymbol{\epsilon}}^{\alpha}$ as in Section 7.1, the only difference being that the expressions for the phase-specific chemical potentials become somewhat more complex due to the additional - but known as a function of \mathbf{c}^{α} - elastic contributions.

More precisely, with

$$\begin{aligned} \frac{\partial(\boldsymbol{\mu}_{el}^\alpha \cdot \boldsymbol{\zeta}^\alpha)}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha &= \frac{\partial}{\partial \mathbf{c}^\alpha} \left(-\boldsymbol{\sigma}^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}_{el}^\alpha) : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) + \frac{1}{2} \left(\boldsymbol{\epsilon}_{el}^\alpha : \left(\frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) : \boldsymbol{\epsilon}_{el}^\alpha \right) \right) \cdot \boldsymbol{\chi}^\alpha \\ &= -\boldsymbol{\sigma}^\alpha : \left(\frac{\partial^2 \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial (\mathbf{c}^\alpha)^2} : (\boldsymbol{\zeta}^\alpha \otimes \boldsymbol{\chi}^\alpha) \right) - \left(\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha \right) \cdot \frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \\ &\quad + \frac{1}{2} \boldsymbol{\epsilon}_{el}^\alpha : \left(\frac{\partial^2 \mathbf{C}^\alpha}{\partial (\mathbf{c}^\alpha)^2} : (\boldsymbol{\zeta}^\alpha \otimes \boldsymbol{\chi}^\alpha) \right) : \boldsymbol{\epsilon}_{el}^\alpha - \left(\left(\frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) : \boldsymbol{\epsilon}_{el}^\alpha \right) : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha \right), \end{aligned}$$

or, using

$$\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \mathbf{c}^\alpha} \cdot d\mathbf{c}^\alpha = -\mathbf{C}^\alpha : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot d\mathbf{c}^\alpha \right) + \left(\frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}^\alpha} \cdot d\mathbf{c}^\alpha \right) : \boldsymbol{\epsilon}_{el}^\alpha \quad (7.192)$$

for replacing the derivative of the stress w.r.t. \mathbf{c}^α , one obtains the formula

$$\begin{aligned} \frac{\partial^2 f_{el}^\alpha}{\partial (\mathbf{c}^\alpha)^2} : (\boldsymbol{\zeta}^\alpha \otimes \boldsymbol{\chi}^\alpha) &= \frac{1}{2} \boldsymbol{\epsilon}_{el}^\alpha : \left(\frac{\partial^2 \mathbf{C}^\alpha}{\partial (\mathbf{c}^\alpha)^2} : (\boldsymbol{\zeta}^\alpha \otimes \boldsymbol{\chi}^\alpha) \right) : \boldsymbol{\epsilon}_{el}^\alpha - \boldsymbol{\sigma}^\alpha : \left(\frac{\partial^2 \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial (\mathbf{c}^\alpha)^2} : (\boldsymbol{\zeta}^\alpha \otimes \boldsymbol{\chi}^\alpha) \right) \\ &\quad - \left(\left(\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha \right) : \boldsymbol{\epsilon}_{el}^\alpha \right) : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) - \left(\left(\frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) : \boldsymbol{\epsilon}_{el}^\alpha \right) : \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha \right) \quad (7.193) \\ &\quad + \left(\frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\chi}^\alpha \right) : \mathbf{C}^\alpha : \left(\frac{\partial \boldsymbol{\epsilon}_{el}^\alpha}{\partial \mathbf{c}^\alpha} \cdot \boldsymbol{\zeta}^\alpha \right) \end{aligned}$$

for characterizing the second derivative of the mechanical contribution to the free energy density with respect to the phase-specific concentration.

Once the \mathbf{c}^α have been determined, the evaluation of the stress $\boldsymbol{\sigma}$ is then straightforward.

Remark 161. This can nevertheless exacerbate the issues already discussed in Remark 104 when considering the \mathbf{c}^α as functions of the chemical potential, since any additional contribution further decreases the likelihood of being able to perform the conversion from $(\boldsymbol{\mu}, \boldsymbol{\epsilon})$ to \mathbf{c}^α explicitly. \diamond

In terms of the driving force $\frac{\partial f}{\partial \phi^\alpha}$ for the phasefield equation, it is clear that, due to the choice $\boldsymbol{\epsilon}^\beta = \boldsymbol{\epsilon}$, this approach entails very little changes with respect to the previous considerations and one has

$$\frac{\partial f}{\partial \phi^\alpha} = \sum_{\beta=1}^N \left(f_{ch}^\beta(\mathbf{c}^\alpha) + \frac{1}{2} (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\beta) : \mathbf{C}^\beta : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^\beta) - \boldsymbol{\mu} \cdot \mathbf{c}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha}.$$

The Reuss-Sachs Case

In contrast (unless is no dependence on \mathbf{c}^α in f_{el}^α), using a Reuss-Sachs-type approach already entails some notable complications, since the defining requirement $\boldsymbol{\sigma}^\alpha = \boldsymbol{\sigma}$ for some $\boldsymbol{\sigma}$ on the phase-specific stresses $\boldsymbol{\sigma}^\alpha$ is now fixed in terms of system depending on **two** primal unknowns $(\boldsymbol{\epsilon}^\alpha, \mathbf{c}^\alpha)$ and one dual one in terms of $\boldsymbol{\sigma}$,

$$\boldsymbol{\sigma}^\alpha = \mathbf{C}^\alpha(\mathbf{c}^\alpha) : (\boldsymbol{\epsilon}^\alpha - \tilde{\boldsymbol{\epsilon}}^\alpha(\mathbf{c}^\alpha)) \stackrel{!}{=} \boldsymbol{\sigma}. \quad (7.194)$$

One can therefore not proceed in a sequential fashion as in the Voigt-Taylor case and instead has to consider the full coupled chemo-mechanical system simultaneously.

Considering for simplicity an a priori reduced formulation and reusing the notation \mathcal{P}_p for the indices of those phases with $h^\alpha(\phi) > 0$, the basic structure of a Newton-step on the full coupled system is of the form

$$\left(\begin{array}{c} \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \mathbf{c}^\alpha} \right)_{\alpha \in \mathcal{P}_p} & \begin{pmatrix} -\mathbf{I} \\ \vdots \\ -\mathbf{I} \end{pmatrix} \\ \left(h^1 \mathbf{I} \quad \dots \quad h^N \mathbf{I} \right) & \mathbf{0} \end{array} \right) & \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \boldsymbol{\epsilon}^\alpha} \right)_{\alpha \in \mathcal{P}_p} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right) \\ \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \boldsymbol{\epsilon}^\alpha} \right)_{\alpha \in \mathcal{P}_p} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right) & \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \mathbf{c}^\alpha} \right)_{\alpha \in \mathcal{P}_p} & \begin{pmatrix} -\mathbf{I} \\ \vdots \\ -\mathbf{I} \end{pmatrix} \\ \left(h^1 \mathbf{I} \quad \dots \quad h^N \mathbf{I} \right) & \mathbf{0} \end{array} \right) \end{array} \right) \left(\begin{array}{c} \left(\begin{array}{c} (\delta \tilde{\mathbf{c}}^\alpha)_{\alpha \in \mathcal{P}_p} \\ (\delta \boldsymbol{\epsilon}^\alpha)_{\alpha \in \mathcal{P}_p} \\ \boldsymbol{\sigma} \end{array} \right) \right) = \left(\begin{array}{c} \left(\begin{array}{c} (-\tilde{\boldsymbol{\mu}}^\alpha)_{\alpha \in \mathcal{P}_p} \\ \mathbf{r}_{\tilde{\boldsymbol{\epsilon}}} \\ (-\boldsymbol{\sigma}^\alpha)_{\alpha \in \mathcal{P}_p} \\ \mathbf{r}_\boldsymbol{\epsilon} \end{array} \right) \end{array} \right).$$

i.e. a system consisting of two generalized saddle-point problems⁹⁸ coupled by two very sparse off-diagonal blocks.

In addition to the part $\frac{\partial^2 f^\alpha}{\partial(\mathbf{c}^\alpha)^2}$ expressed in Equation (7.193), from Equation (7.191), a variation $d\epsilon^\alpha$ of the phase-specific strain leads to the variation

$$\frac{\partial \boldsymbol{\mu}_{el}^\alpha}{\partial \epsilon^\alpha} : d\epsilon^\alpha = -(\mathbf{C}^\alpha(\mathbf{c}^\alpha) : d\epsilon^\alpha) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^\alpha}{\partial \mathbf{c}^\alpha} + d\epsilon^\alpha : \left(\epsilon_{el}^\alpha : \frac{\partial \mathbf{C}^\alpha}{\partial \mathbf{c}^\alpha} \right) = d\epsilon^\alpha : \frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \mathbf{c}^\alpha} \quad (7.195)$$

of the elastic contribution to the chemical potential, whereas, by the phase-specific stress-strain relationship, one trivially has $\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \epsilon^\alpha} = \mathbf{C}^\alpha$.

As both diagonal subproblems are themselves highly sparse with a known block-diagonal substructure in the “primal” part, it is natural to try to make use of this knowledge by using a block-factorization based approach. There are three relatively obvious choices for doing so.

The first two are to start either analogously to the purely chemical or the purely mechanical case⁹⁹, i.e. by performing a block-elimination step either in terms of the inverse of the full chemical or mechanical saddle-point problems. The principle advantage of this approach is that, from an implementation point of view, this semi-sequential treatment allows maintaining at least some of the modularity from the simpler setting without an explicit coupling between the chemical and mechanical parts. Its major disadvantage is that since both system containing a coupling between all phases due to the sum-constraint on the $\tilde{\mathbf{c}}^\alpha$ resp. the ϵ^α , the first elimination leads to a loss of the remaining block-diagonal structure in the respective other subblock. More precisely, eliminating $((\epsilon^\alpha)_{\alpha \mathcal{P}_p}, \boldsymbol{\sigma})$ in terms of $((\tilde{\mathbf{c}}^\alpha)_{\alpha \mathcal{P}_p}, \tilde{\boldsymbol{\mu}})$ leads in particular to the subblock $\text{diag} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} \right)_{\alpha \in \mathcal{P}_p}$ being replaced by a fully filled matrix, whereas eliminating $((\tilde{\mathbf{c}}^\alpha)_{\alpha \mathcal{P}_p}, \tilde{\boldsymbol{\mu}})$ in terms of $((\epsilon^\alpha)_{\alpha \mathcal{P}_p}, \boldsymbol{\sigma})$ has the same effect on the previously block-diagonal part $\text{diag} \left(\frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \epsilon^\alpha} \right)_{\alpha \in \mathcal{P}_p}$.

This is avoided by the third approach, which consists in eliminating all primal unknowns in terms of the dual ones, i.e. by performing a block-elimination based on the subblocks

$$\frac{\partial^2 \tilde{f}^\alpha}{\partial(\tilde{\mathbf{c}}^\alpha, \epsilon^\alpha)^2} \begin{pmatrix} \delta \tilde{\mathbf{c}}^\alpha \\ \delta \epsilon^\alpha \end{pmatrix} = \begin{pmatrix} \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} & \frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \epsilon^\alpha} \\ \frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \tilde{\mathbf{c}}^\alpha} & \frac{\partial \boldsymbol{\sigma}^\alpha}{\partial \epsilon^\alpha} \end{pmatrix} \begin{pmatrix} \delta \tilde{\mathbf{c}}^\alpha \\ \delta \epsilon^\alpha \end{pmatrix} = \begin{pmatrix} \tilde{\boldsymbol{\mu}} - \tilde{\boldsymbol{\mu}}^\alpha \\ \boldsymbol{\sigma} - \boldsymbol{\sigma}^\alpha \end{pmatrix}$$

for each phase for determining $(\tilde{\mathbf{c}}^\alpha, \epsilon^\alpha)$ as functions of $\tilde{\boldsymbol{\mu}}$ and $\boldsymbol{\sigma}$, and thus reducing the problem to the one of solving a coupled Schur-complement equation for $(\tilde{\boldsymbol{\mu}}, \boldsymbol{\sigma})$. A clear advantage of such a factorization is that it makes the highest use of the underlying sparsity pattern. A potential drawback from an implementation point of view is that it is the “most explicitly coupled” of the three approaches above, and therefore, while likely the most efficient choice, the one which is the trickiest with respect to code-reusability.

In combination with the Reuss-Sachs model, the dependence of both \mathbf{c}^β and ϵ^β on ϕ leads to the same two contributions as before, i.e. $-\boldsymbol{\mu} \cdot \mathbf{c}^\beta$ from the chemical part and $-\boldsymbol{\sigma} : \epsilon^\beta$ from the mechanical part,

$$\frac{\partial f}{\partial \phi^\alpha} = \sum_{\beta=1}^N \left(f_{ch}^\beta(\mathbf{c}^\alpha) + \frac{1}{2}(\epsilon - \tilde{\epsilon}^\beta) : \mathbf{C}^\beta : (\epsilon - \tilde{\epsilon}^\beta) - \boldsymbol{\mu} \cdot \mathbf{c}^\beta - \boldsymbol{\sigma} : \epsilon^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha}.$$

Remark 162. From similar calculations as before, it also follows that the important properties $\boldsymbol{\mu} = \frac{\partial f}{\partial \mathbf{c}}$ and $\boldsymbol{\sigma} = \frac{\partial f}{\partial \epsilon}$ still hold with $\boldsymbol{\sigma}$ and $\boldsymbol{\mu}$ as defined above, despite the dependence of the

⁹⁸The more standard saddle-point structure could again be recovered by “undoing” the division by h^α in the upper lines of each of the subblocks.

⁹⁹The procedure outlined above essentially being a nonlinear (in the \mathbf{c}^α) variation of the latter approach which does, after performing the linearization in the $\delta \mathbf{c}^\alpha$, reduce to the same resulting set of equations.

phase-specific quantities on both \mathbf{c} and $\boldsymbol{\epsilon}$. For the chemical potential, this follows from

$$\frac{\partial f}{\partial \mathbf{c}} = \sum_{\alpha} \left(\frac{\partial f^{\alpha}}{\partial \mathbf{c}^{\alpha}} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{c}} + \frac{\partial f^{\alpha}}{\partial \boldsymbol{\epsilon}^{\alpha}} : \frac{\partial \boldsymbol{\epsilon}^{\alpha}}{\partial \mathbf{c}} \right) h^{\alpha}(\phi) = \sum_{\alpha} \left((\boldsymbol{\mu} + \lambda^{\alpha} \mathbf{e}) \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{c}} + \boldsymbol{\sigma} : \frac{\partial \boldsymbol{\epsilon}^{\alpha}}{\partial \mathbf{c}} \right) h^{\alpha}(\phi)$$

combined with $\mathbf{e} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{c}} = \frac{\partial(\mathbf{e} \cdot \mathbf{c}^{\alpha})}{\partial \mathbf{c}} = \frac{\partial 1}{\partial \mathbf{c}} = 0$ as well as, after extracting the phase-independent terms out of the sums, $\sum_{\alpha} \frac{\partial \mathbf{c}^{\alpha}}{\partial \mathbf{c}} h^{\alpha} = \frac{\partial \mathbf{c}}{\partial \mathbf{c}}$ and $\sum_{\alpha} \frac{\partial \boldsymbol{\epsilon}^{\alpha}}{\partial \mathbf{c}} h^{\alpha} = \frac{\partial \boldsymbol{\epsilon}}{\partial \mathbf{c}} = \mathbf{0}$. The calculation for the stresses is completely analogous. \diamond

Remark 163. As the “mechanical” subsystem is linear in both the $\boldsymbol{\epsilon}^{\alpha}$ and $\boldsymbol{\sigma}$, yet another alternative for solving this system is to proceed in the same manner as for the original description of the Reuss-Sachs model and eliminate the the $\boldsymbol{\epsilon}^{\alpha}$ as functions of the \mathbf{c}^{α} and $\boldsymbol{\sigma}$ **before** performing the linearization. This first leads to

$$\boldsymbol{\epsilon}^{\alpha}(\mathbf{c}^{\alpha}, \boldsymbol{\sigma}) = \mathbf{S}^{\alpha}(\mathbf{c}^{\alpha}) : \boldsymbol{\sigma} + \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha}),$$

and, combining this with the sum-constraint on the $\boldsymbol{\epsilon}^{\alpha}$, to the Schur-complement system

$$\boldsymbol{\epsilon} = \sum_{\alpha=1}^N \boldsymbol{\epsilon}^{\alpha}(\mathbf{c}^{\alpha}, \boldsymbol{\sigma}) h^{\alpha}(\phi) = \mathbf{S}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}) : \boldsymbol{\sigma} + \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N})$$

for $\boldsymbol{\sigma}$, where $\tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}) = \sum_{\alpha} \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha}) h^{\alpha}(\phi)$. It follows that

$$\boldsymbol{\sigma}^{\alpha}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) = \boldsymbol{\sigma}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) = \mathbf{C}_{RS}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}) : \left(\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}) \right)$$

and thus

$$\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha}) = \boldsymbol{\epsilon}_{el}^{\alpha}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) = \mathbf{S}^{\alpha}(\mathbf{c}^{\alpha}) : \boldsymbol{\sigma}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}),$$

the only difference to the purely mechanical case being that $\boldsymbol{\sigma}$ now in addition depends on all the $(\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}$. Reinserting these expressions into Equation (7.191) leads to

$$\begin{aligned} \boldsymbol{\mu}^{\alpha} &= \boldsymbol{\mu}_{ch}^{\alpha}(\mathbf{c}^{\alpha}) - \left(\mathbf{C}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) \right) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^{\alpha}}{\partial \mathbf{c}^{\alpha}} + \frac{1}{2} (\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \left((\boldsymbol{\epsilon}^{\alpha} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) : \frac{\partial \mathbf{C}^{\alpha}}{\partial \mathbf{c}^{\alpha}} \right) \\ &= \boldsymbol{\mu}_{ch}^{\alpha}(\mathbf{c}^{\alpha}) - \boldsymbol{\sigma}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) : \frac{\partial \tilde{\boldsymbol{\epsilon}}^{\alpha}}{\partial \mathbf{c}^{\alpha}} \\ &\quad + \frac{1}{2} \left(\mathbf{S}^{\alpha}(\mathbf{c}^{\alpha}) : \boldsymbol{\sigma}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) \right) : \left(\left(\mathbf{S}^{\alpha}(\mathbf{c}^{\alpha}) : \boldsymbol{\sigma}(\boldsymbol{\phi}, (\mathbf{c}^{\alpha})_{1 \leq \alpha \leq N}, \boldsymbol{\epsilon}) \right) : \frac{\partial \mathbf{C}^{\alpha}}{\partial \mathbf{c}^{\alpha}} \right), \end{aligned}$$

and thus a system in the \mathbf{c}^{α} alone.

A subsequent linearization for dealing with the nonlinearity in \mathbf{c}^{α} will then leads to essentially the same system as obtained above by linearization first and then performing a block-elimination of the mechanical quantities in terms of $\tilde{\mathbf{c}}^{\alpha}$ and $\tilde{\boldsymbol{\mu}}$. \diamond

The Jump-Condition Based Formulation

Finally, if there are notable differences in the mechanical behavior of the different phases, it may be beneficial to complement the conditions on the chemical potentials with a mechanical model as in Section 7.2.2 or 7.2.4.

The Two-Phase Case In the two-phase case, the simplest choice is to, as in Equation (7.88), impose $\epsilon^1(\boldsymbol{\epsilon}, \mathbf{a}) = \boldsymbol{\epsilon} - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S$ and $\epsilon^2(\boldsymbol{\epsilon}, \mathbf{a}) = \boldsymbol{\epsilon} + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S$, and to fix \mathbf{a} through the additional condition

$$\mathbf{r}_\sigma := \llbracket \boldsymbol{\sigma}^\alpha(\mathbf{c}^\alpha, \boldsymbol{\epsilon}^\alpha) \rrbracket^{12} \cdot \mathbf{n} = (\boldsymbol{\sigma}^2(\mathbf{c}^2, \boldsymbol{\epsilon}^2) - \boldsymbol{\sigma}^1(\mathbf{c}^1, \boldsymbol{\epsilon}^1)) \cdot \mathbf{n} \stackrel{!}{=} \mathbf{0}. \quad (7.196)$$

Since the $\boldsymbol{\epsilon}^\alpha$ are both given in terms of a simple explicit expression in terms of $\boldsymbol{\epsilon}$ and \mathbf{a} , it is clearly convenient to rewrite the system consisting of Equation (7.196) and the equality of the (reduced) chemical potentials as functions of $(\mathbf{c}^\alpha, \tilde{\boldsymbol{\mu}})$ and \mathbf{a} instead of the $\boldsymbol{\epsilon}^\alpha$ themselves. A Newton-step on this system consists in solving

$$\left(\begin{array}{c} \left(\begin{array}{ccc} \frac{\partial \tilde{\boldsymbol{\mu}}^1}{\partial \tilde{\mathbf{c}}^1} & \mathbf{0} & -\mathbf{I} \\ \mathbf{0} & \frac{\partial \tilde{\boldsymbol{\mu}}^2}{\partial \tilde{\mathbf{c}}^2} & -\mathbf{I} \\ h^1 \mathbf{I} & h^2 \mathbf{I} & \mathbf{0} \end{array} \right) & \left(\begin{array}{c} \frac{\partial \tilde{\boldsymbol{\mu}}^1}{\partial \mathbf{a}} \\ \frac{\partial \tilde{\boldsymbol{\mu}}^2}{\partial \mathbf{a}} \\ \mathbf{0} \end{array} \right) \\ \left(\begin{array}{ccc} \frac{\partial \mathbf{r}_\sigma}{\partial \tilde{\mathbf{c}}^1} & \frac{\partial \mathbf{r}_\sigma}{\partial \tilde{\mathbf{c}}^2} & \mathbf{0} \end{array} \right) & \frac{\partial \mathbf{r}_\sigma}{\partial \mathbf{a}} \end{array} \right) \left(\begin{array}{c} \delta \tilde{\mathbf{c}}^1 \\ \delta \tilde{\mathbf{c}}^2 \\ \tilde{\boldsymbol{\mu}} \\ \delta \mathbf{a} \end{array} \right) = \left(\begin{array}{c} -\tilde{\boldsymbol{\mu}}^1 \\ -\tilde{\boldsymbol{\mu}}^2 \\ \mathbf{r}_\tilde{\mathbf{c}} \\ \mathbf{r}_\sigma \end{array} \right), \quad (7.197)$$

and thus a relatively simple addition to the purely chemical problem¹⁰⁰. In addition, from Equation (7.195) together with the expressions for ϵ^1 and ϵ^2 , it follows immediately that the $\frac{\partial \tilde{\boldsymbol{\mu}}^\alpha}{\partial \mathbf{a}}$ are characterized by

$$\frac{\partial \tilde{\boldsymbol{\mu}}^1}{\partial \mathbf{a}} = -h^2 \frac{\partial(\boldsymbol{\sigma}^1(\tilde{\mathbf{c}}^1, \mathbf{a}) \cdot \mathbf{n})}{\partial \tilde{\mathbf{c}}^1} \quad \text{and} \quad \frac{\partial \tilde{\boldsymbol{\mu}}^2}{\partial \mathbf{a}} = h^1 \frac{\partial(\boldsymbol{\sigma}^2(\tilde{\mathbf{c}}^2, \mathbf{a}) \cdot \mathbf{n})}{\partial \tilde{\mathbf{c}}^2},$$

whereas the remaining new entries are given by

$$\frac{\partial \mathbf{r}_\sigma}{\partial \tilde{\mathbf{c}}^1} = -\frac{\partial(\boldsymbol{\sigma}^1(\mathbf{c}^1, \mathbf{a}) \cdot \mathbf{n})}{\partial \tilde{\mathbf{c}}^1} \quad \text{and} \quad \frac{\partial \mathbf{r}_\sigma}{\partial \tilde{\mathbf{c}}^2} = +\frac{\partial(\boldsymbol{\sigma}^2(\mathbf{c}^2, \mathbf{a}) \cdot \mathbf{n})}{\partial \tilde{\mathbf{c}}^2}$$

as well as the matrix

$$\frac{\partial \mathbf{r}_\sigma}{\partial \mathbf{a}} = \mathbf{n} \cdot \left(h^2(\phi) \mathbf{C}^1(\mathbf{c}^1) + h^1(\phi) \mathbf{C}^2(\mathbf{c}^2) \right) \cdot \mathbf{n}$$

already used in Section 7.2.2 in a non-incremental form in \mathbf{a} .

Remark 164. It can be noted that this matrix is, after “undoing” the various divisions by the h -functions (i.e. after a multiplication of the first two rows with h^1 and h^2 respectively and the fourth one by $h^1 h^2$ symmetric, as is to be expected from the variational characterization underlying the equality of the chemical potentials as in [25] (see Section 7.1) and the continuity condition on the stresses as proposed in [51] (see Section 7.2.2). \diamond

Even though the coupling of the local chemical and mechanical equilibrium conditions now implicitly defines \mathbf{c}^α and the $\boldsymbol{\epsilon}^\alpha$ as functions of the four given parameters $(\phi, \mathbf{c}, \boldsymbol{\epsilon}, \mathbf{n})$, the usual calculations show that this does not really affect the previous contributions in the separate chemical and mechanical case. Firstly, for the driving force, one has

$$\frac{\partial f}{\partial \phi^\alpha} = \sum_\beta f^\beta(\mathbf{c}^\beta, \boldsymbol{\epsilon}^\beta) \frac{\partial h^\beta}{\partial \phi^\alpha} + \sum_\beta \frac{\partial f^\beta}{\partial \mathbf{c}^\beta} \cdot \frac{\partial \mathbf{c}^\beta}{\partial \phi^\alpha} + \frac{\partial f}{\partial \boldsymbol{\epsilon}^\beta} : \frac{\partial \boldsymbol{\epsilon}^\beta}{\partial \phi^\alpha} h^\beta.$$

The summation over the first term in the second sum leads as previously to $\sum_\beta -\boldsymbol{\mu} \cdot \mathbf{c}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$. The last term consists as in Section 7.2.2 of two contributions, one given by the explicit dependence of the definitions of $\epsilon^1 - h^2(\phi)(\mathbf{a} \otimes \mathbf{n})_S$ and $\epsilon^2 + h^1(\phi)(\mathbf{a} \otimes \mathbf{n})_S$ on ϕ and the other one due to

¹⁰⁰Some points regarding potential solution strategies will be discussed in relation with the multiphase formulation.

the implicit dependence of \mathbf{a} on ϕ . The latter one cancels as the vanishing of the normal stress jump ensures that $\sum_{\beta} \frac{\partial f^{\beta}}{\partial \mathbf{a}} h^{\beta} = \mathbf{0}$, whereas the first one results in

$$h^1 \boldsymbol{\sigma}^1 : (\mathbf{a} \otimes \mathbf{n})_S \left(-\frac{\partial h^2}{\partial \phi^{\alpha}}\right) + h^2 \boldsymbol{\sigma}^2 : (\mathbf{a} \otimes \mathbf{n})_S \frac{\partial h^1}{\partial \phi^{\alpha}}.$$

Using $(\mathbf{a} \otimes \mathbf{n})_S = \boldsymbol{\epsilon}^2 - \boldsymbol{\epsilon}^1$, $h^1 + h^2 = 0$ and $\frac{\partial h^1}{\partial \phi^{\alpha}} = -\frac{\partial h^2}{\partial \phi^{\alpha}}$, this expression can, as in Section 7.2.2, be rewritten into a more pleasant form as $\sum_{\beta} -\boldsymbol{\sigma} : \boldsymbol{\epsilon}^{\beta} \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}$, leaving the final expression

$$\frac{\partial f}{\partial \phi^{\alpha}} = \sum_{\beta} \left(f^{\beta}(\mathbf{c}^{\beta}, \boldsymbol{\epsilon}^{\beta}) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}} - \boldsymbol{\mu} \cdot \mathbf{c}^{\beta} - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^{\beta} \right) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}.$$

The remaining contribution to the phasefield equation through $\frac{\partial f}{\partial \mathbf{n}}$ can also directly be taken from Equation (7.99), since, even though the \mathbf{c}^{α} now also depend on \mathbf{n} , one has

$$\sum_{\beta} \frac{\partial f^{\beta}}{\partial \mathbf{c}^{\beta}} \cdot \frac{\partial \mathbf{c}^{\beta}}{\partial \mathbf{n}} h^{\beta} = \sum_{\beta} (\boldsymbol{\mu} + \lambda^{\beta} \mathbf{e}) \cdot \frac{\partial \mathbf{c}^{\beta}}{\partial \mathbf{n}} h^{\beta} = \boldsymbol{\mu} \cdot \sum_{\beta} \frac{\partial \mathbf{c}^{\beta}}{\partial \mathbf{n}} h^{\beta} = \boldsymbol{\mu} \cdot \frac{\partial (\sum_{\beta} \mathbf{c}^{\beta} h^{\beta})}{\partial \mathbf{n}} = \boldsymbol{\mu} \cdot \frac{\partial \mathbf{c}}{\partial \mathbf{n}} = \mathbf{0}$$

Finally, making use of the analogous arguments (and in particular $\frac{\partial f}{\partial \mathbf{a}} = \mathbf{0}$), one can further verify that $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ continue to satisfy $\boldsymbol{\mu} = \frac{\partial f}{\partial \mathbf{c}}$ and $\boldsymbol{\sigma} = \frac{\partial f}{\partial \boldsymbol{\epsilon}}$.

The Multiphase Case Extending the two-phase model above to a multiphase setting is again a highly non-trivial task as it inherits the same difficulties already encountered in the purely mechanical setting in Section 7.2.4 and will only be discussed for the two models in [62] and [63]. The simplest to extend is the one based on a common normal vector as the model was seen to essentially preserve the pleasant mathematical properties from the two-phase case. Using the jump-vector based reformulation with $\boldsymbol{\epsilon}^{\alpha} = \boldsymbol{\epsilon} - (\mathbf{a}^{\alpha} \otimes \mathbf{n})_S$, the system to be solved is obtained by complementing the chemical conditions with the equality of the phase-specific normal stresses

$$\boldsymbol{\sigma}^{\alpha}(\mathbf{c}^{\alpha}, \boldsymbol{\epsilon}^{\alpha}) \cdot \mathbf{n} = \mathbf{n} \cdot \left(\mathbf{C}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon} - (\mathbf{a}^{\alpha} \otimes \mathbf{n})_S - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) \right) \stackrel{!}{=} \boldsymbol{\sigma} \cdot \mathbf{n}$$

with the yet to be determined total normal stress $\boldsymbol{\sigma} \cdot \mathbf{n}$. For the application of the basic Newton-scheme, it is again convenient to explicitly make use of the expressions for the $\boldsymbol{\epsilon}^{\alpha}$ in terms of the jump-vectors. In the simpler reduced form, this leads to the system

$$\left(\begin{array}{c} \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^{\alpha}}{\partial \tilde{\mathbf{c}}^{\alpha}} \right)_{\alpha \in \mathcal{P}_p} & \begin{pmatrix} -I \\ \vdots \\ -I \end{pmatrix} \\ \left(h^1 \mathbf{I} \quad \dots \quad h^N \mathbf{I} \right) & \mathbf{0} \end{array} \right) & \left(\begin{array}{cc} \text{diag} \left(\frac{\partial \tilde{\boldsymbol{\mu}}^{\alpha}}{\partial \mathbf{a}^{\alpha}} \right)_{\alpha \in \mathcal{P}_p} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right) \\ \left(\begin{array}{cc} \text{diag} \left(\frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n})}{\partial \tilde{\mathbf{c}}^{\alpha}} \right)_{\alpha \in \mathcal{P}_p} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right) & \left(\begin{array}{cc} \text{diag} \left(\frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n})}{\partial \mathbf{a}^{\alpha}} \right)_{\alpha \in \mathcal{P}_p} & \begin{pmatrix} -I \\ \vdots \\ -I \end{pmatrix} \\ \left(h^1 \mathbf{I} \quad \dots \quad h^N \mathbf{I} \right) & \mathbf{0} \end{array} \right) \end{array} \right) \left(\begin{array}{c} \left(\begin{array}{c} (\delta \tilde{\mathbf{c}}^{\alpha})_{\alpha \in \mathcal{P}_p} \\ (\delta \mathbf{a}^{\alpha})_{\alpha \in \mathcal{P}_p} \\ \boldsymbol{\sigma} \cdot \mathbf{n} \end{array} \right) \end{array} \right) = \left(\begin{array}{c} \left(\begin{array}{c} (-\tilde{\boldsymbol{\mu}}^{\alpha})_{\alpha \in \mathcal{P}_p} \\ (-\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n})_{\alpha \in \mathcal{P}_p} \\ \mathbf{r}_{\mathbf{a}} \end{array} \right) \end{array} \right),$$

where, from Equation (7.195) and $d\boldsymbol{\epsilon}^{\alpha} = -(\mathbf{d}\mathbf{a}^{\alpha} \otimes \mathbf{n})_S$, one has $\frac{\partial \tilde{\boldsymbol{\mu}}^{\alpha}}{\partial \mathbf{a}^{\alpha}} = -\frac{\partial \boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}}{\partial \tilde{\mathbf{c}}^{\alpha}}$ and $\mathbf{r}_{\boldsymbol{\sigma}^{\alpha}} := \boldsymbol{\sigma} \cdot \mathbf{n} - \boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}$, while $\mathbf{r}_{\mathbf{a}} = \sum_{\alpha=1}^N \mathbf{a}^{\alpha} h^{\alpha}(\phi)$ is the error in the sum-constraint on the average of the jump-vectors.

Remark 165. This system clearly has the same basic structure as the one in the Reuss-Sachs case with the \mathbf{a}^{α} replacing the $\boldsymbol{\epsilon}^{\alpha}$ and in particular the sum-constraint $\sum_{\alpha} \mathbf{a}^{\alpha} h^{\alpha} = \mathbf{0}$ replacing $\sum_{\alpha} \boldsymbol{\epsilon}^{\alpha} h^{\alpha} = \boldsymbol{\epsilon}$. Even though it is in principle somewhat easier to solve due to the smaller number of unknowns in the \mathbf{a}^{α} and $\boldsymbol{\sigma} \cdot \mathbf{n}$ as compared to the full values of $\boldsymbol{\epsilon}^{\alpha}$ and $\boldsymbol{\sigma}$, the considerations underlying its solution are essentially the same as for the Reuss-Sachs based model. \diamond

The driving force follows by combining the standard argument for the chemical part with the arguments preceding Equations (7.170) and (7.171)¹⁰¹ to be given by

$$\frac{\partial}{\partial \phi^\alpha} \Big|_{\mathbf{n}} f(\phi, \boldsymbol{\epsilon}) = \sum_{\beta} (f_{el}^{\beta} - \boldsymbol{\mu} \cdot \mathbf{c}^{\beta} + (\boldsymbol{\sigma}^{\beta} \cdot \mathbf{n}) \cdot \mathbf{a}^{\beta}) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}, \quad (7.198)$$

or, using $\boldsymbol{\epsilon}^{\beta} = \boldsymbol{\epsilon} - (\mathbf{a}^{\beta} \otimes \mathbf{n})_S$, by

$$\frac{\partial f_{el}}{\partial \phi^{\alpha}} \Big|_{\mathbf{n}} = \sum_{\beta=1}^N (f_{el}^{\beta}(\boldsymbol{\epsilon}^{\beta}) - \boldsymbol{\mu} \cdot \mathbf{c}^{\beta} - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^{\beta}) \frac{\partial h^{\beta}}{\partial \phi^{\alpha}}. \quad (7.199)$$

The differentiation with respect to \mathbf{n} leads to the same result as in Equation (7.172), since, despite the additional dependence of the \mathbf{c}^{β} on \mathbf{n} , their total contributions drop out again due to the relation satisfied with respect to the chemical potential. Finally, $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ as obtained in the solution of the local quasi-equilibrium conditions above still satisfy, despite the additional implicit dependencies, $\boldsymbol{\mu} = \frac{\partial f}{\partial \mathbf{c}}$ and $\boldsymbol{\sigma} = \frac{\partial f}{\partial \boldsymbol{\epsilon}}$ ¹⁰².

The natural extension of the model by Schneider [61] is obtained by simply admitting an additional dependence on the \mathbf{c}^{α} in the normal continuity conditions

$$\llbracket \boldsymbol{\sigma}^{\alpha} \rrbracket^{\text{R}\alpha} \cdot \mathbf{n}^{\text{R}\alpha} = (\boldsymbol{\sigma}^{\alpha}(\mathbf{c}^{\alpha}, \boldsymbol{\epsilon}^{\alpha}) - \boldsymbol{\sigma}^{\text{R}}(\mathbf{c}^{\text{R}}, \boldsymbol{\epsilon}^1)) \cdot \mathbf{n}^{\text{R}\alpha} \stackrel{!}{=} \mathbf{0}$$

with respect to the reference phase combined with the previous definition of the phase-specific strains $\boldsymbol{\epsilon}^{\text{R}} = \boldsymbol{\epsilon} - \sum_{\beta \neq \text{R}} (\mathbf{a}^{\text{R}\beta} \otimes \mathbf{n}^{\text{R}\beta})_S h^{\beta}(\phi)$ and $\boldsymbol{\epsilon}^{\alpha} = \boldsymbol{\epsilon}^{\text{R}} + (\mathbf{a}^{\text{R}\alpha} \otimes \mathbf{n}^{\text{R}\alpha})_S$, $\alpha \neq \text{R}$ in Equation (7.174). In terms of a Newton-step for the \mathbf{c}^{α} and $\mathbf{a}^{\text{R}\alpha}$, one has, from Equation (7.195), to consider the additional increments

$$\sum_{\beta \neq \text{R}} \frac{\partial \boldsymbol{\mu}_{el}^1}{\partial \mathbf{a}^{\text{R}\beta}} \cdot \delta \mathbf{a}^{\text{R}\beta} = - \sum_{\beta \neq \text{R}} \delta \mathbf{a}^{\text{R}\beta} \cdot h^{\beta} \frac{\partial (\boldsymbol{\sigma}^{\text{R}} \cdot \mathbf{n}^{\text{R}\beta})}{\partial \mathbf{c}^{\text{R}}} \quad (7.200)$$

in the chemical potential for the reference phase, and

$$\begin{aligned} \sum_{\beta \neq \text{R}} \frac{\partial \boldsymbol{\mu}_{el}^{\alpha}}{\partial \mathbf{a}^{\text{R}\beta}} \cdot \delta \mathbf{a}^{\text{R}\beta} &= - \sum_{\beta \neq \text{R}} \delta \mathbf{a}^{\text{R}\beta} \cdot h^{\beta} \frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}^{\text{R}\beta})}{\partial \mathbf{c}^{\alpha}} + \delta \mathbf{a}^{\text{R}\alpha} \cdot \frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}^{\text{R}\alpha})}{\partial \mathbf{c}^{\alpha}} \\ &= \delta \mathbf{a}^{\text{R}\alpha} \cdot (1 - h^{\alpha}) \frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}^{\text{R}\alpha})}{\partial \mathbf{c}^{\alpha}} - \sum_{\beta \neq \text{R}, \alpha} \delta \mathbf{a}^{\text{R}\beta} \cdot h^{\beta} \frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}^{\text{R}\beta})}{\partial \mathbf{c}^{\alpha}} \end{aligned} \quad (7.201)$$

for $\alpha \neq \text{R}$, as well as the increments

$$\frac{\partial (\boldsymbol{\sigma}^{\alpha} \cdot \mathbf{n}^{\text{R}\alpha})}{\partial \mathbf{c}^{\alpha}} \cdot \delta \mathbf{c}^{\alpha} - \frac{\partial (\boldsymbol{\sigma}^{\text{R}} \cdot \mathbf{n}^{\text{R}\alpha})}{\partial \mathbf{c}^{\text{R}}} \cdot \delta \mathbf{c}^{\text{R}} \quad (7.202)$$

in the normal stress jumps due to the changes in the concentration, whereas the expression for $\sum_{\beta \neq \text{R}} \frac{\partial (\boldsymbol{\sigma}^{\alpha} - \boldsymbol{\sigma}^{\text{R}}) \cdot \mathbf{n}^{\text{R}\alpha}}{\partial \mathbf{a}^{\text{R}\beta}} \cdot \delta \mathbf{a}^{\text{R}\beta}$ is, by linearity in the $\mathbf{a}^{\text{R}\beta}$ and Equation (7.176) given by

$$\left(\mathbf{n}^{\text{R}\alpha} \cdot ((1 - h^{\alpha}) \mathbf{c}^{\alpha} + h^{\alpha} \mathbf{c}^{\text{R}}) \cdot \mathbf{n}^{\text{R}\alpha} \right) \cdot \delta \mathbf{a}^{\text{R}\alpha} - \sum_{\beta \neq \text{R}, \alpha} h^{\beta} \left(\mathbf{n}^{\text{R}\alpha} \cdot (\mathbf{c}^{\alpha} - \mathbf{c}^{\text{R}}) \cdot \mathbf{n}^{\text{R}\beta} \right) \cdot \delta \mathbf{a}^{\text{R}\beta}. \quad (7.203)$$

¹⁰¹I.e. essentially a combination of the common normal stress for all phases and the constraint equation $\sum_{\alpha} \mathbf{a}^{\alpha} h^{\alpha} = \mathbf{0}$.

¹⁰²For example for $\boldsymbol{\mu}$, the potentially problematic new contribution could arise from the dependency of \mathbf{a}^{α} on \mathbf{c} . Since $\sum_{\alpha} \boldsymbol{\sigma}^{\alpha} : \left(\frac{\partial \mathbf{a}^{\alpha}}{\partial c_i} \otimes \mathbf{n} \right) h^{\alpha} = (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \sum_{\alpha} \frac{\partial \mathbf{a}^{\alpha}}{\partial c_i} h^{\alpha} = \frac{\partial \sum_{\alpha} \mathbf{a}^{\alpha} h^{\alpha}}{\partial c_i} = \mathbf{0}$ by the constraint $\sum_{\alpha} \mathbf{a}^{\alpha} h^{\alpha} = \mathbf{0}$, this dependency drops out.

The total system therefore has a quite different structure from the one based on using a common normal vector. In particular, it consists of a fully filled (1,2) and (2,2)-block by Equations (7.201) and (7.203). For this reason, the only sparsity pattern which can be made use of lies in the chemical subblock, i.e. two reasonable alternatives to solving the full system is to either perform a block-elimination on the concentration increments in terms of the jump-vectors and the chemical potential or a full block-elimination on the chemical quantities in terms of the jump-vectors.

The driving force calculation is, except for the standard additional term $-\boldsymbol{\mu} \cdot \mathbf{c}^\beta$ essentially the same as the one for Equations (7.170) and (7.171), leading to

$$\frac{\partial}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} f(\boldsymbol{\phi}, \boldsymbol{\epsilon}) = \sum_{\beta} \left(f^\beta - \boldsymbol{\mu} \cdot \mathbf{c}^\beta - \boldsymbol{\sigma} : \boldsymbol{\epsilon}^\beta \right) \frac{\partial h^\beta}{\partial \phi^\alpha} + \sum_{\beta \neq R} h^\beta \left((\boldsymbol{\sigma}^\beta - \boldsymbol{\sigma}) \cdot \mathbf{n}^{R\beta} \right) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}}. \quad (7.204)$$

As before, there is a remaining dependence on the $\frac{\partial \mathbf{a}^{R\beta}}{\partial \phi^\alpha}$. For a full evaluation, one therefore has to obtain the $\frac{\partial \mathbf{a}^{R\beta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}}$ from a linearization of the local equilibrium conditions¹⁰³. The non-variational structure of the equilibrium conditions with respect to the \mathbf{a} is also the reason why, similarly to the relation $\boldsymbol{\sigma} = \frac{\partial f_{el}}{\partial \boldsymbol{\epsilon}}$ not holding within multiphase regions, the ‘‘chemical potential’’ obtained from the solution of the local conditions above does not actually satisfy $\boldsymbol{\mu} = \frac{\partial f}{\partial \mathbf{c}}$ ¹⁰⁴. As both points lead to notable complications and there is already a variational consistency in terms of $\boldsymbol{\sigma}$ and $\frac{\partial f}{\partial \boldsymbol{\epsilon}}$, the pragmatic choice is again to ignore this difficulty within the multiphase regions¹⁰⁵.

¹⁰³At least assuming that the reference phase does not change as there will usually not be differentiability at these points. Differentiating the equilibrium conditions $\tilde{\boldsymbol{\mu}}^\beta = \tilde{\boldsymbol{\mu}}$ for all β together with the sum-constraint on the $\tilde{\mathbf{c}}^\beta$ and $(\boldsymbol{\sigma}^\alpha - \boldsymbol{\sigma}^R) \cdot \mathbf{n}^{R\beta} = \mathbf{0}$ for all $\beta \neq R$ with respect to ϕ^α with the normal vectors kept fixed one obtains

$$\frac{\partial \tilde{\boldsymbol{\mu}}^\beta}{\partial \tilde{\mathbf{c}}^\beta} \frac{\partial \tilde{\mathbf{c}}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} + \frac{\partial \tilde{\boldsymbol{\mu}}^\beta}{\partial \mathbf{a}^{R\delta}} \frac{\partial \mathbf{a}^{R\delta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} - \frac{\partial \tilde{\boldsymbol{\mu}}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} = \frac{\partial \tilde{\boldsymbol{\mu}}^\beta}{\partial \boldsymbol{\epsilon}^\beta} : \sum_{\delta \neq R} (\mathbf{a}^{R\delta} \otimes \mathbf{n}^{R\delta})_S \frac{\partial h^\delta}{\partial \phi^\alpha}$$

as well as $\sum_{\beta} h^\beta \frac{\partial \tilde{\mathbf{c}}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} = -\sum_{\beta} \mathbf{c}^\beta \frac{\partial h^\beta}{\partial \phi^\alpha}$ from the chemical part which needs to be combined with the conditions

$$\begin{aligned} & -\frac{\partial \boldsymbol{\sigma}^R \cdot \mathbf{n}^{R\beta}}{\partial \tilde{\mathbf{c}}^R} \frac{\partial \tilde{\mathbf{c}}^R}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} + \frac{\partial \boldsymbol{\sigma}^\beta \cdot \mathbf{n}^{R\beta}}{\partial \tilde{\mathbf{c}}^\beta} \frac{\partial \tilde{\mathbf{c}}^\beta}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} + \sum_{\delta \neq R} \frac{\partial (\boldsymbol{\sigma}^\beta - \boldsymbol{\sigma}^R) \cdot \mathbf{n}^{R\beta}}{\partial \mathbf{a}^{R\delta}} \frac{\partial \mathbf{a}^{R\delta}}{\partial \phi^\alpha} \Big|_{\mathbf{n}^{\alpha\beta}} \\ & = -\left(\frac{\partial \boldsymbol{\sigma}^\beta \cdot \mathbf{n}^{R\beta}}{\partial \boldsymbol{\epsilon}^\beta} - \frac{\partial \boldsymbol{\sigma}^R \cdot \mathbf{n}^{R\beta}}{\partial \boldsymbol{\epsilon}^R} \right) : \sum_{\delta \neq R} (\mathbf{a}^{R\delta} \otimes \mathbf{n}^{R\delta})_S \frac{\partial h^\delta}{\partial \phi^\alpha} \end{aligned}$$

for all $\beta \neq R$.

If not based on an interpolation function satisfying $\frac{\partial h^\delta}{\partial \phi^\alpha} = 0$ if $h^\delta(\boldsymbol{\phi}) = 0$, there is a slight additional complication here as compared to the original solution procedure for the local equilibrium conditions, since, even though the coupling in the unknowns is still only through those phases with $h^\beta(\boldsymbol{\phi}) > 0$, this is not true for the right-hand sides.

¹⁰⁴Instead, one has

$$\frac{\partial f}{\partial \mathbf{c}} = \sum_{\alpha} (\boldsymbol{\mu} + \lambda^\alpha \mathbf{e}) \frac{\partial \mathbf{c}^\alpha}{\partial \mathbf{c}} h^\alpha + \sum_{\alpha} \boldsymbol{\sigma}^\alpha : \frac{\partial \boldsymbol{\epsilon}^\alpha}{\partial \mathbf{c}} h^\alpha.$$

While the first term again simplifies to $\boldsymbol{\mu}$, the second one maintains a dependence of the jump vectors on the total concentration through

$$\begin{aligned} & -\sum_{\alpha \neq R} \sum_{\beta \neq R} h^\alpha h^\beta (\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}^{R\beta}) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \mathbf{c}} + \sum_{\alpha \neq R} h^\alpha (\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}^{R\alpha}) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \mathbf{c}} = -\sum_{\beta \neq R} h^\beta (\boldsymbol{\sigma} \cdot \mathbf{n}^{R\beta}) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \mathbf{c}} + \sum_{\alpha \neq R} h^\alpha (\boldsymbol{\sigma}^\alpha \cdot \mathbf{n}^{R\alpha}) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \mathbf{c}} \\ & = \sum_{\alpha \neq R} h^\alpha \left((\boldsymbol{\sigma}^\alpha - \boldsymbol{\sigma}) \cdot \mathbf{n}^{R\alpha} \right) \cdot \frac{\partial \mathbf{a}^{R\beta}}{\partial \mathbf{c}}. \end{aligned}$$

¹⁰⁵Recall that within two-phase regions, there is no such issue.

Remark 166. One can in principle also construct a chemo-mechanical extension of the model by Tschukin by e.g. simply adding a concentration dependence in his expression

$$\frac{1}{2} \sum_{\alpha=1}^N \left(\boldsymbol{\Sigma}^{\alpha}(\mathbf{c}^{\alpha}) : (\boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}^{\alpha}(\mathbf{c}^{\alpha})) - \frac{1}{2} \sum_{\beta \neq \alpha} \llbracket \boldsymbol{\Sigma} \rrbracket^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) : \bar{\boldsymbol{\mathcal{S}}}_{nn}^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) : \llbracket \boldsymbol{\Sigma} \rrbracket^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) h^{\beta}(\boldsymbol{\phi}) \right) h^{\alpha}(\boldsymbol{\phi})$$

for the mechanical free energy density¹⁰⁶ with the \mathbf{c}^{α} determined based on the usual chemical equilibrium condition and then defining $\boldsymbol{\sigma} = \frac{\partial f}{\partial \boldsymbol{\epsilon}}$ as in [74].

As one of the motivations of the expressions in [74] seems to have been to obtain a fairly explicit expression of the underlying free energy in terms of the parameters $(\boldsymbol{\phi}, \nabla \boldsymbol{\phi}, \boldsymbol{\epsilon})$, this is somewhat contrary in spirit as the \mathbf{c}^{α} would then be defined implicitly as $\mathbf{c}^{\alpha} = \mathbf{c}^{\alpha}(\boldsymbol{\phi}, \nabla \boldsymbol{\phi}, \mathbf{c}, \boldsymbol{\epsilon})$ and has not been investigated in more detail¹⁰⁷. \diamond

¹⁰⁶As discussed previously, the term in the double-summation could also be rewritten as

$$\left(\llbracket \boldsymbol{\Sigma} \rrbracket^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) \cdot \mathbf{n}^{\alpha\beta} \right) \cdot \left(\mathbf{n}^{\alpha\beta} \cdot \bar{\boldsymbol{\mathcal{C}}}^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) \cdot \mathbf{n}^{\alpha\beta} \right)^{-1} \cdot \left(\llbracket \boldsymbol{\Sigma} \rrbracket^{\alpha\beta}(\mathbf{c}^{\alpha}, \mathbf{c}^{\beta}) \cdot \mathbf{n}^{\alpha\beta} \right)$$

for avoiding the use of the pseudo-inverse $\bar{\boldsymbol{\mathcal{S}}}_{nn}^{\alpha\beta}$.

¹⁰⁷It may be interesting to do so though. In this regard, it should be noted that this implicit dependence is, again by the variational characterization of the \mathbf{c}^{α} , not really problematic here. In particular, the basic form of the expression for the final stress in Equation (7.181) would be maintained, since, despite the a priori additional appearance of the terms $\frac{\partial f}{\partial \mathbf{c}^{\alpha}} \cdot \frac{\partial \mathbf{c}^{\alpha}}{\partial \boldsymbol{\epsilon}}$, these would again drop out in the sum making use of the equality of the chemical potentials and $\sum_{\alpha} \mathbf{c}^{\alpha} h^{\alpha} = 1$. For similar reasons, most of the description in Section 7.2.4 is expected to carry over without major modifications.

Chapter 8

Summary

The focus of this thesis has been on some primarily practical implications of different constraints arising in phasefield based models.

Chapter 6 discussed several consequences of the Gibbs-Simplex constraint on the phasefield values. Even though the constraint itself is relatively simple, it has far-reaching numerical consequences. One of these was studied in detail in terms of the spatial discretization for a simple but prototypical two-phase problem, illustrating some interesting relations between the numerical interface width as compared to the parameter ϵ and the associated discrete energetics. The central point of Section 6.3 is, besides a short summary of some by now mostly standard background such as the LROP approach from [40], the description of a simple local projection algorithm in combination with the mobility-matrix based dynamics as in [68]. Even though similar algorithms are likely known in the quadratic programming literature, the author is not aware of any specific source with an explicit description adapted to the particular nature of \mathbf{M}^1 .

In contrast, Chapter 7 considers coupled problems with “internal” phase-specific variables in a chemical, mechanical and chemo-mechanical setting. Subsections 7.1.2 and 7.1.3 are essentially introductory in nature, adding some additional details to the very elegant description of the quantitative free energy model as described in [25]. This is then compared in Subsection 7.1.4 with the now more commonly used description in [56] and [19]. This discussion is in a sense continued in Sections 7.1.5 and 7.1.5, where some practical implications of the choice of “primary” unknown in terms of the concentration or the chemical potential are discussed. In particular, it was argued that, while the latter description can be highly efficient for particularly simple free energy densities, the former approach should be preferred in the general case. Finally, Subsection 7.1.6 considered an extension of the model to a non-isothermal situation, illustrating, among other things, the pleasant interplay of the variational definitions on two classical thermodynamic relations.

A closely related modeling approach in the mechanical case is discussed in Section 7.2. As the conditions fixing the phase-specific quantities is more complex in this case, various different formulations have been suggested in the literature even in the two-phase case. Some advantages and disadvantages of these formulations as well as their links are summarized in Subsections 7.2.2 and 7.2.3 as a basis for the analysis of some of the difficulties encountered by three extensions to the multiphase setting by Schneider and Tschukin in [63], [62] and [74]. As a final topic, Subsection 7.2.5 outlined several practical and theoretical implications obtained by coupling the chemical model from Section 7.1 with the various mechanical models.

¹There are also surprisingly few explicit descriptions of the Euclidian projection onto the Gibbs-simplex as many algorithms are instead designed to handle weighted projections subject to general equality and inequality constraints. Nevertheless, it is clear that the weight being given by the identity, the equality constraint by $\mathbf{e} \cdot \phi = 1$ and the inequality constraints by $\phi^\alpha \geq 0$ allows for some significant simplifications as compared to the general case.

The central aspect of Chapter 7 is not the descriptions of the various models themselves, which are for the most part either well-known or relatively straightforward extensions of already known approaches, but to sensitize to and analyze some of the issues which can arise due to the introduction of more complex models, in particular in the presence of additional “internal” unknowns in terms of phasespecific quantities. Since the advantages of such models for the approximation quality of the related sharp interface modes even in the presence of artificially large interfaces are by now clearly recognized, it is to be expected that their use will become even more widespread in the years to come. As seen on the chemical and mechanical examples considered in Sections 7.1 and 7.2, there are some pitfalls when trying to combine such approaches with the very common two-step formalism for the derivation of new phasefield models in the literature. On the one hand, the discussion of the multiphase mechanical model in [63] illustrated that independently postulating a phasefield functional and a set of equations to be satisfied by an additional set of unknowns besides the phasefield itself entails the risk of introducing an incompatibility with the standard variational approach. Even though this is not necessarily a problem in itself as the resulting model may still deliver very accurate results, one has to be careful about invoking an intuitively very pleasant but ultimately nonexistent energy minimization principle for its justification.

On the other hand, Chapter 7 also highlighted potential fallacies arising from the derivation of the drivingforces for the phasefield equation through a purely formal partial differentiation of a given functional with respect to ϕ , in particular also in combination with potential changes of primary unknowns. The questions raised by an interdependence between various unknowns and equations are of course well-understood and have mostly been answered a long time ago. Unfortunately, the fact that this allow to justify such a procedure for many of the earlier and simpler models, this seems to have lead to a fairly widespread misconception in the more applied phasefield community that this is simply “what needs to be done” in order to minimize or maximize the underlying functional. One of the purposes of the discussion of the models is therefore also to stress the importance of underlying variational principles or the lack thereof and the additional drivingforce contributions arising due to both local constraints imposed on the phasespecific quantities and global constraints imposed e.g. through a global mass-conservation principle. It is the author’s hope that the detailed consideration of some of these issues in the very applied context of Chapter 7 may help to clear up some a priori confusing but ultimately clear-cut questions, in particular for young researches newly discovering the fascinating domain of phasefield modeling.

Appendix A

Calculation of the Discrete One-Dimensional Phasefield Energy and Interface Width

This section will provide some more details on the calculations underlying the one-dimensional discrete profile in subsection 6.2.3, and in particular the expressions (repeated here for convenience)

$$\mathcal{E}(N) = \frac{4\gamma}{\pi^2\epsilon}N + \frac{\epsilon\gamma}{(\Delta x)^2} \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa) = \frac{4\gamma}{\pi^2\epsilon} \left(N + \cot\left(\frac{\kappa N}{2}\right) \sqrt{\frac{\pi^2\epsilon^2}{4(\Delta x)^2} - 1} \right) \quad (\text{A.1})$$

for the total energy,

$$\sum_{i=1}^{N-1} \frac{1}{\epsilon} w(\phi_i) = \frac{8}{\pi^2\epsilon} \gamma \left(\frac{1}{4} \left(1 - \cot^2\left(\frac{\kappa N}{2}\right) \right) N + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \cot(\kappa) \right) \quad (\text{A.2})$$

for the total contribution by the bulk potential in Equation (6.42) and

$$\sum_{i=1}^{N-1} \epsilon a_i = \frac{8}{\pi^2\epsilon} \gamma \left(\frac{1}{4} \left(1 + \cot^2\left(\frac{\kappa N}{2}\right) \right) N + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \frac{1}{\sin(\kappa)} \right), \quad (\text{A.3})$$

for the total contribution by the gradient energy in Equation (6.43) as well as the actual number of interface points from Equation (6.47),

$$N = N_{min} = \left\lceil \frac{2}{\kappa} \tan^{-1} \left(\sqrt{\frac{\pi^2\epsilon^2}{4(\Delta x)^2} - 1} \right) \right\rceil. \quad (\text{A.4})$$

A.1 A Quick Recap of the First-Order Analysis

As seen in Subsection 6.2.3, fixing for convenience $\phi_0 = 0$ as the last bulk-point on the left and $\phi_N = 1$ as the first bulk-point on the right, the discrete first-order optimality condition (6.31) **within** the interface (i.e. at all points $1 \leq i \leq N-1$ for which the strict inequality $0 < \phi_i < 1$ holds) together with imposing the conditions $\phi_0 = 0$ and $\phi_N = 1$ at the transition to the bulk implies

$$\phi_i = \frac{1}{2} - \frac{1}{2} \cos(\kappa i) + \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa i) \quad (\text{A.5})$$

with κ determined by

$$\cos(\kappa) = 1 - \frac{8(\Delta x)^2}{\pi^2 \epsilon^2} = 1 - \frac{1}{2} \left(\frac{4\Delta x}{\pi \epsilon} \right)^2 \quad (\text{A.6})$$

$$\sin(\kappa) = \sqrt{1 - \left(1 - \frac{8(\Delta x)^2}{\pi^2 \epsilon^2}\right)^2} = \sqrt{2 \frac{8(\Delta x)^2}{\pi^2 \epsilon^2} - \left(\frac{8(\Delta x)^2}{\pi^2 \epsilon^2}\right)^2} = \frac{8(\Delta x)^2}{\pi^2 \epsilon^2} \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}. \quad (\text{A.7})$$

Remark 167. Note that up to this point, this is simply a consequence of the linear difference equation satisfied within the interface and the choice of an (in total) increasing profile going from 0 on the left to 1 on the right and will up to this point work for an arbitrary integer $N > 0$. In order to select the correct N , one still needs take both the underlying demand of a minimization of the total energy $\Delta x \mathcal{E}(N) = \sum_{i=1}^{N-1} (\epsilon a_i + \frac{1}{\epsilon} w_i) \Delta x$ and the constraints $0 \leq \phi_i \leq 1$ on the profile into account. \diamond

The positivity constraint $\phi_1 \geq 0$ leads, within the relevant range¹, to the restriction

$$N \leq N_{max} = \left\lceil \frac{\pi}{\kappa} \right\rceil \quad (\text{A.8})$$

on the total number of points, whereas the sign-constraint on the multipliers μ^\pm for the box-constraint $0 \leq \phi_i \leq 1$ lead to the lower bound

$$N \geq N_{min} := \left\lceil \frac{\pi}{\kappa} - 1 \right\rceil \quad (\text{A.9})$$

Comparing the bounds it is obvious that this leaves two potential choices for N . Determining the optimal profile is therefore not possible based on the first-order necessary condition alone and requires a more detailed energetic analysis, which is of an inherent interest anyway.

A.2 The Discrete Energetic Analysis

Determining the total energy starting from a given sinusoidal profile is essentially a matter of a simple integration. The only difficulty is that, unlike in the continuous case, the correct interpretation of this integration here is actually a summation over a set of discrete sinusoidal values. While this generally leads to - provided one is even able to do this analytically - very cumbersome expressions, it can be achieved in a very elegant fashion in the solenoidal case by the use of **Lagrange's trigonometric identities**

$$\sum_{i=1}^{N-1} \sin(\lambda i) = \frac{1}{2} \cot\left(\frac{\lambda}{2}\right) - \frac{\cos\left(\lambda\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\lambda}{2}\right)}, \quad (\text{A.10})$$

$$\sum_{i=1}^{N-1} \cos(\lambda i) = -\frac{1}{2} + \frac{\sin\left(\lambda\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\lambda}{2}\right)}. \quad (\text{A.11})$$

In fact, a simple differentiation (based on a finite-dimensional vector of ϕ -values) shows that the first-order optimality condition as in Equation (6.38) corresponds to the minimization of the discrete energy function (the fixed factor Δx essentially being irrelevant for the minimizer, but of course important from a physical point of view)

$$\begin{aligned} \Delta x \mathcal{E}(N) &= \Delta x \epsilon \sum_{i=0}^N \frac{1}{2} \left(\left(\frac{\phi_i - \phi_{i-1}}{\Delta x} \right)^2 + \left(\frac{\phi_{i+1} - \phi_i}{\Delta x} \right)^2 \right) + \frac{16\Delta x}{\pi^2 \epsilon} \sum_{i=0}^N \phi_i (1 - \phi_i) \\ &= \epsilon \Delta x \sum_{i=1}^N \left(\frac{\phi_i - \phi_{i-1}}{\Delta x} \right)^2 + \frac{16\Delta x}{\pi^2 \epsilon} \sum_{i=1}^{N-1} \phi_i (1 - \phi_i) \end{aligned}$$

¹By periodicity, significantly larger choices of N would in principle also be possible, but would be related to transitions involving several sinusoidals, which is obviously energetically unfavorable.

where the contributions from all $i < 0$ and $i > N$ vanish due to both the zero gradients and the vanishing of the w -term at $\phi \in \{0, 1\}$.

Based on the variational character of the phasefield Equation (6.31) within the interface, it turns out that the total energy is significantly easier to evaluate than the bulk potential or gradient energy contributions themselves. In fact, since

$$\begin{aligned} \sum_{i=1}^N \left(\frac{\phi_i - \phi_{i-1}}{\Delta x} \right)^2 &= \frac{1}{\Delta x} \left(\sum_{i=1}^N \frac{\phi_i - \phi_{i-1}}{\Delta x} \phi_i - \sum_{i=0}^{N-1} \frac{\phi_{i+1} - \phi_i}{\Delta x} \phi_i \right) \\ &= \frac{1}{\Delta x} \left(\sum_{i=1}^{N-1} \frac{\phi_i - \phi_{i-1}}{\Delta x} \phi_i - \sum_{i=1}^{N-1} \frac{\phi_{i+1} - \phi_i}{\Delta x} \phi_i \right) + \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_1 \right) \\ &= - \sum_{i=1}^{N-1} \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} \phi_i + \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right), \end{aligned}$$

the total energy can equivalently be rewritten as

$$\begin{aligned} \mathcal{E}(N) &= \epsilon \left(- \sum_{i=1}^{N-1} \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} \phi_i + \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right) \right) + \frac{16}{\pi^2 \epsilon} \sum_{i=1}^{N-1} \phi_i (1 - \phi_i) \\ &= \sum_{i=1}^{N-1} \left(- \epsilon \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} + \frac{16}{\pi^2 \epsilon} (1 - \phi_i) \right) \phi_i + \epsilon \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right) \\ &= \frac{1}{2} \sum_{i=1}^{N-1} \underbrace{\left(- 2\epsilon \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} + \frac{16}{\pi^2 \epsilon} (1 - 2\phi_i) \right)}_{=0 \text{ for } 1 \leq i \leq N-1} \phi_i \\ &\quad + \frac{1}{2} \sum_{i=1}^{N-1} \frac{16}{\pi^2 \epsilon} \phi_i + \epsilon \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right) \\ &= \sum_{i=1}^{N-1} \frac{8}{\pi^2 \epsilon} \phi_i + \epsilon \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right). \end{aligned} \tag{A.12}$$

With $\phi_i = \frac{1}{2} - \frac{1}{2} \cos(\kappa i) + c_2 \sin(\kappa i)$ where

$$c_2 := \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \tag{A.13}$$

and Lagrange's trigonometric identities in Equation (A.10), this can be "integrated" due to

$$\sum_{i=1}^{N-1} \left(\frac{1}{2} - \frac{1}{2} \cos(\kappa i) + c_2 \sin(\kappa i) \right) = \frac{N-1}{2} - \frac{1}{2} \left(-\frac{1}{2} + \frac{\sin\left(\kappa\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\kappa}{2}\right)} \right) + c_2 \left(\frac{1}{2} \cot\left(\frac{\kappa}{2}\right) - \frac{\cos\left(\kappa\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\kappa}{2}\right)} \right)$$

and therefore leads to

$$\begin{aligned} \mathcal{E}(N) &= \frac{8}{\pi^2 \epsilon} \left[\frac{N-1}{2} - \frac{1}{2} \left(-\frac{1}{2} + \frac{\sin\left(\kappa\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\kappa}{2}\right)} \right) + c_2 \left(\frac{1}{2} \cot\left(\frac{\kappa}{2}\right) - \frac{\cos\left(\kappa\left(N - \frac{1}{2}\right)\right)}{2 \sin\left(\frac{\kappa}{2}\right)} \right) \right] \\ &\quad + \epsilon \frac{1}{\Delta x} \left(\frac{\phi_N - \phi_{N-1}}{\Delta x} \phi_N - \frac{\phi_1 - \phi_0}{\Delta x} \phi_0 \right). \end{aligned}$$

Combined with the choice of orientation for the interface here, i.e. $\phi_0 = 0$ and $\phi_N = 1$, one further has

$$\phi_N - \phi_{N-1} = (1 - \cos(\kappa)) \left(\frac{1}{2} \right) + \sin(\kappa) \left(\frac{1}{2} \sin(\kappa N) + c_2 \cos(\kappa N) \right)$$

and thus

$$\mathcal{E} = \frac{8}{\pi^2 \epsilon} \left[\frac{N-1}{2} - \frac{1}{2} \left(-\frac{1}{2} + \frac{\sin(\kappa(N-\frac{1}{2}))}{2 \sin(\frac{\kappa}{2})} \right) + c_2 \left(\frac{1}{2} \cot(\frac{\kappa}{2}) - \frac{\cos(\kappa(N-\frac{1}{2}))}{2 \sin(\frac{\kappa}{2})} \right) \right] + \frac{\epsilon}{(\Delta x)^2} \left(\frac{1}{2} (1 - \cos(\kappa)) + \sin(\kappa) \left(\frac{1}{2} \sin(\kappa N) + c_2 \cos(\kappa N) \right) \right).$$

Based on the half-angle formula

$$\cot\left(\frac{\theta}{2}\right) = \frac{1 + \cos(\theta)}{\sin(\theta)} \quad (\text{A.14})$$

for the cosine, it follows that $c_2 = \frac{1 + \cos(\kappa N)}{2 \sin(\kappa N)}$, where κ is independent of N . Combining this with the addition formulae

$$\sin(\kappa(2N-1)) = \sin(2\kappa N) \cos(\kappa) - \cos(2\kappa N) \sin(\kappa), \quad (\text{A.15})$$

$$\cos(\kappa(2N-1)) = \cos(2\kappa N) \cos(\kappa) + \sin(2\kappa N) \sin(\kappa), \quad (\text{A.16})$$

for the sine and cosine, one obtains

$$\mathcal{E} = \frac{8}{\pi^2 \epsilon} \left[\frac{N-1}{2} - \frac{1}{2} \left(-\frac{1}{2} + \frac{1}{2} \left(\sin(\kappa N) \cot(\frac{\kappa}{2}) - \cos(\kappa N) \right) \right) + c_2 \left(\frac{1}{2} \cot(\frac{\kappa}{2}) - \frac{1}{2} \left(\cos(\kappa N) \cot(\frac{\kappa}{2}) + \sin(\kappa N) \right) \right) \right] + \frac{\epsilon}{(\Delta x)^2} \left(\frac{1}{2} (1 - \cos(\kappa)) - \sin(\kappa) \left(\frac{1}{2} \sin(\kappa N) + c_2 \cos(\kappa N) \right) \right).$$

A significant simplification can be obtained due to $c_2 \sin(\kappa N) = \frac{1}{2} (1 + \cos(\kappa N))$ and therefore,

$$c_2 \cos(\kappa N) = \frac{\cos(\kappa N) + (1 - \sin^2(\kappa N))}{2 \sin(\kappa N)} = c_2 - \frac{1}{2} \sin(\kappa N)$$

from which it follows that

$$\begin{aligned} & -\frac{1}{2} \left(-\frac{1}{2} + \frac{1}{2} \left(\sin(\kappa N) \cot(\frac{\kappa}{2}) - \cos(\kappa N) \right) \right) + c_2 \left(\frac{1}{2} \cot(\frac{\kappa}{2}) - \frac{1}{2} \left(\cos(\kappa N) \cot(\frac{\kappa}{2}) + \sin(\kappa N) \right) \right) \\ &= \frac{1}{4} - \frac{1}{4} (\sin(\kappa N) + 2c_2 \cos(\kappa N)) \cot(\frac{\kappa}{2}) + \frac{1}{4} \cos(\kappa N) + \frac{1}{2} c_2 \cot(\frac{\kappa}{2}) - \frac{1}{2} c_2 \sin(\kappa N) \\ &= \frac{1}{4} - \frac{1}{4} (\sin(\kappa N) + 2(c_2 - \frac{1}{2} \sin(\kappa N))) \cot(\frac{\kappa}{2}) + \frac{1}{4} \cos(\kappa N) + \frac{1}{2} c_2 \cot(\frac{\kappa}{2}) - \frac{1}{4} (1 + \cos(\kappa N)) = 0. \end{aligned}$$

Further making use of

$$\frac{1}{2} \sin(\kappa N) + c_2 \cos(\kappa N) = \frac{1}{2} \sin(\kappa N) + c_2 - \frac{1}{2} \sin(\kappa N) = c_2$$

in the last term, the remaining expression

$$\mathcal{E} = \frac{4}{\pi^2 \epsilon} (N-1) + \frac{\epsilon}{(\Delta x)^2} \left(\frac{1}{2} (1 - \cos(\kappa)) + c_2 \sin(\kappa) \right)$$

then solely depends on the ‘‘width’’ of the interface and an additional term arising at the transition to the bulk. As $\frac{\epsilon}{(\Delta x)^2} (1 - \cos(\kappa)) = \frac{8}{\pi^2 \epsilon}$, this reduces to

$$\mathcal{E} = \frac{4}{\pi^2 \epsilon} N + \frac{\epsilon}{(\Delta x)^2} c_2 \sin(\kappa) = \frac{4}{\pi^2 \epsilon} N + \frac{\epsilon}{(\Delta x)^2} \frac{1}{2} \cot\left(\frac{\kappa N}{2}\right) \sin(\kappa), \quad (\text{A.17})$$

or, inserting the expression $\sin(\kappa) = \frac{8(\Delta x)^2}{\pi^2 \epsilon^2} \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}$ from Equation (A.7), to

$$\mathcal{E} = \frac{4}{\pi^2 \epsilon} N + c_2 \frac{8}{\pi^2 \epsilon} \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1} = \frac{4}{\pi^2 \epsilon} N + \frac{4}{\pi^2 \epsilon} \cot\left(\frac{\kappa N}{2}\right) \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}.$$

With $\cot\left(\frac{\kappa}{2}\right) = \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}$, one then finally obtains to the expression

$$\mathcal{E} = \frac{4}{\pi^2 \epsilon} \left(N + 2c_2 \cot\left(\frac{\kappa}{2}\right) \right) = \frac{4}{\pi^2 \epsilon} \left(N + \cot\left(\frac{\kappa N}{2}\right) \cot\left(\frac{\kappa}{2}\right) \right) \quad (\text{A.18})$$

in Equation (6.46).

The evaluation of the individual contributions through the gradient- and bulk energy contributions is unfortunately a little more tedious than that of the total energy as one cannot directly make use of the equation itself to eliminate the quadratic terms in ϕ as in Equation (A.12). Starting with the somewhat simpler bulk-potential term, one has to evaluate

$$\begin{aligned} \Delta x \sum \frac{1}{\epsilon} w_i &= \frac{16\Delta \bar{x}}{\pi^2} \sum_{i=1}^{N-1} \left(\frac{1}{2} - \frac{1}{2} \cos(\kappa i) + c_2 \sin(\kappa i) \right) \left(\frac{1}{2} - \left(-\frac{1}{2} \cos(\kappa i) + c_2 \sin(\kappa i) \right) \right) \\ &= \frac{16\Delta \bar{x}}{\pi^2} \sum_{i=1}^{N-1} \frac{1}{4} - \left(\frac{1}{2} \cos(\kappa i) - c_2 \sin(\kappa i) \right)^2 \\ &= \frac{16\Delta \bar{x}}{\pi^2} \sum_{i=1}^{N-1} \frac{1}{4} - \left(\frac{1}{4} \cos^2(\kappa i) - c_2 \sin(\kappa i) \cos(\kappa i) + c_2^2 \sin^2(\kappa i) \right). \end{aligned}$$

The quadratic terms in the sinusoidal function can be eliminated using the basic identities $\sin(\theta) \cos(\theta) = \frac{1}{2} \sin(2\theta)$, $\sin^2(\theta) = \frac{1 - \cos(2\theta)}{2}$ and $\cos^2(\theta) = \frac{1 + \cos(2\theta)}{2}$, leading to

$$\begin{aligned} \Delta x \sum \frac{1}{\epsilon} w_i &= \frac{16\Delta \bar{x}}{\pi^2} \sum_{i=1}^{N-1} \frac{1}{4} - \left(\frac{1}{4} \frac{1 + \cos(2\kappa i)}{2} - \frac{1}{2} c_2 \sin(2\kappa i) + c_2^2 \frac{1 - \cos(2\kappa i)}{2} \right) \\ &= \frac{8\Delta \bar{x}}{\pi^2 \epsilon} \sum_{i=1}^{N-1} \left(\frac{1}{4} - c_2^2 \right) - \left(\frac{1}{4} - c_2^2 \right) \cos(2\kappa i) + c_2 \sin(2\kappa i). \end{aligned}$$

This is now in a form where Equation (A.10) can be applied, resulting (with $\lambda = 2\kappa$) in

$$\Delta x \sum \frac{1}{\epsilon} w_i = \frac{8\Delta \bar{x}}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) - \left(\frac{1}{4} - c_2^2 \right) \left(-\frac{1}{2} + \frac{\sin(\kappa(2N-1))}{2 \sin(\kappa)} \right) \right) + c_2 \left(\frac{1}{2} \cot(\kappa) - \frac{\cos(\kappa(2N-1))}{2 \sin(\kappa)} \right).$$

Using the addition formulae (A.15) this can further be expanded to

$$\begin{aligned} \Delta x \sum \frac{1}{\epsilon} w_i &= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) + \left(c_2^2 - \frac{1}{4} \right) \left(-\frac{1}{2} + \frac{1}{2} \sin(2\kappa N) \cot(\kappa) - \frac{1}{2} \cos(2\kappa N) \right) \right. \\ &\quad \left. + c_2 \left(\frac{1}{2} \cot(\kappa) - \frac{1}{2} \cos(2\kappa N) \cot(\kappa) - \frac{1}{2} \sin(2\kappa N) \right) \right), \end{aligned}$$

or, with the double angle fomulae

$$\sin(2\kappa N) = 2 \sin(\kappa N) \cos(\kappa N) \quad \text{and} \quad \cos(2\kappa N) = \cos^2(\kappa N) - \sin^2(\kappa N), \quad (\text{A.19})$$

to

$$\begin{aligned}
\Delta x \sum \frac{1}{\epsilon} w_i &= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) + \frac{1}{2} \left(c_2^2 - \frac{1}{4} \right) \left(-1 + 2 \sin(\kappa N) \cos(\kappa N) \cot(\kappa) - \cos^2(\kappa N) + \sin^2(\kappa N) \right) \right. \\
&\quad \left. + \frac{1}{2} c_2 \left(\cot(\kappa) - (\cos^2(\kappa N) - \sin^2(\kappa N)) \cot(\kappa) - 2 \sin(\kappa N) \cos(\kappa N) \right) \right) \\
&= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) + \left(c_2^2 - \frac{1}{4} \right) \left(\sin(\kappa N) \cos(\kappa N) \cot(\kappa) - \cos^2(\kappa N) \right) \right. \\
&\quad \left. + c_2 \left(\sin^2(\kappa N) \cot(\kappa) - \sin(\kappa N) \cos(\kappa N) \right) \right).
\end{aligned}$$

Making further progress now requires reintroducing the definition $\frac{1}{2} \cot\left(\frac{\kappa N}{2}\right)$ of c_2 . First, the half-angle formula (A.14) allows modifying the last term to

$$c_2 \left(\sin^2(\kappa N) \cot(\kappa) - \sin(\kappa N) \cos(\kappa N) \right) = \frac{1}{2} (1 + \cos(\kappa N)) \left(\sin(\kappa N) \cot(\kappa) - \cos(\kappa N) \right).$$

It further implies on the one hand that

$$c_2 \sin^2(\kappa N) = \frac{1}{2} \sin(\kappa N) + \frac{1}{2} \sin(\kappa N) \cos(\kappa N)$$

and on the other hand that

$$\begin{aligned}
c_2^2 - \frac{1}{4} &= \frac{1}{2} c_2 \frac{1 + \cos(\kappa N)}{\sin(\kappa N)} - \frac{1}{4} = \frac{2c_2 + 2c_2 \cos(\kappa N) - \sin(\kappa N)}{4 \sin(\kappa N)} = \frac{2c_2 + 2(c_2 - \frac{1}{2} \sin(\kappa N)) - \sin(\kappa N)}{4 \sin(\kappa N)} \\
&= \frac{4c_2 - 2 \sin(\kappa N)}{4 \sin(\kappa N)} = \frac{c_2}{\sin(\kappa N)} - \frac{1}{2},
\end{aligned}$$

and thus the two additional relations

$$\begin{aligned}
\left(c_2^2 - \frac{1}{4} \right) \sin(\kappa N) &= c_2 - \frac{1}{2} \sin(\kappa N) \\
\left(c_2^2 - \frac{1}{4} \right) \sin(\kappa N) \cos(\kappa N) &= c_2 \cos(\kappa N) - \frac{1}{2} \sin(\kappa N) \cos(\kappa N).
\end{aligned}$$

Combining these observations, one obtains

$$\begin{aligned}
\Delta x \sum \frac{1}{\epsilon} w_i &= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) + \left(c_2^2 - \frac{1}{4} \right) \left(\sin(\kappa N) \cos(\kappa N) \cot(\kappa) - \cos^2(\kappa N) \right) \right. \\
&\quad \left. + c_2 \left(\sin^2(\kappa N) \cot(\kappa) - \sin(\kappa N) \cos(\kappa N) \right) \right) \\
&= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) - \left(c_2^2 - \frac{1}{4} \right) \cos^2(\kappa N) \right. \\
&\quad \left. + \left(c_2 \cos(\kappa N) - \frac{1}{2} \sin(\kappa N) \cos(\kappa N) \right) \cot(\kappa) \right. \\
&\quad \left. - c_2 \sin(\kappa N) \cos(\kappa N) + \left(\frac{1}{2} \sin(\kappa N) + \frac{1}{2} \sin(\kappa N) \cos(\kappa N) \right) \cot(\kappa) \right) \\
&= \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) (N-1) - \left(c_2^2 - \frac{1}{4} \right) \cos^2(\kappa N) + c_2 \cos(\kappa N) \cot(\kappa) \right. \\
&\quad \left. - c_2 \sin(\kappa N) \cos(\kappa N) + \frac{1}{2} \sin(\kappa N) \cot(\kappa) \right).
\end{aligned}$$

Again based on the half-angle formula (A.14), it further follows that $c_2 \cos(\kappa N) = c_2 - \frac{1}{2} \sin(\kappa N)$ and replacing $\cos^2(\kappa N)$ with $1 - \sin^2(\kappa N)$ leads to the major simplification,

$$\Delta x \sum \frac{1}{\epsilon} w_i = \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) N + \left(c_2^2 - \frac{1}{4} \right) \sin^2(\kappa N) + c_2 \cot(\kappa) - c_2 \sin(\kappa N) \cos(\kappa N) \right),$$

and then finally with

$$\left(c_2^2 - \frac{1}{4}\right) \sin^2(\kappa N) - c_2 \sin(\kappa N) \cos(\kappa N) = \left(c_2 - \frac{1}{2} \sin(\kappa N)\right) \sin(\kappa N) - \left(c_2 - \frac{1}{2} \sin(\kappa N)\right) \sin(\kappa N) = 0$$

to the expression

$$\Delta x \sum \frac{1}{\epsilon} w_i = \frac{8}{\pi^2 \epsilon} \left(\left(\frac{1}{4} - c_2^2 \right) N + c_2 \cot(\kappa) \right) \quad (\text{A.20})$$

for the contributions by the bulk-potential already given in Equation (6.42) resp. (A.2).

The energetic contribution of the gradient energy density term in Equation (A.3) resp. (6.43) can in principle be recovered by a similarly tedious calculation using a slight variation of the arguments above, but, as one already disposes of an expression for the total energy and the contribution by the bulk-potential, simply by taking the difference between the expression in Equation (A.1) and Equation (A.2).

Based on the expression for the total energy in Equation (A.1), it is now an easy matter to determine the value of N with the lowest interface energy. Taking the difference between the energies for two successive values of N shows that

$$\mathcal{E}(N+1) - \mathcal{E}(N) = \frac{4}{\pi^2 \epsilon} + \frac{\epsilon}{(\Delta x)^2} (c_2(N+1) - c_2(N)) \sin(\kappa),$$

where $c_2(N)$ is defined in Equation (A.13). Based on

$$c_2(N+1) - c_2(N) = \frac{1}{2} \left(\cot\left(\frac{\kappa(N+1)}{2}\right) - \cot\left(\frac{\kappa N}{2}\right) \right)$$

and

$$\cot\left(\frac{\kappa(N+1)}{2}\right) = \frac{\cot\left(\frac{\kappa N}{2}\right) \cot\left(\frac{\kappa}{2}\right) - 1}{\cot\left(\frac{\kappa N}{2}\right) + \cot\left(\frac{\kappa}{2}\right)},$$

one has

$$\begin{aligned} 2(c_2(N+1) - c_2(N)) &= \frac{\cot\left(\frac{\kappa N}{2}\right) \cot\left(\frac{\kappa}{2}\right) - 1}{\cot\left(\frac{\kappa N}{2}\right) + \cot\left(\frac{\kappa}{2}\right)} - \cot\left(\frac{\kappa N}{2}\right) = -\frac{1 + \cot^2\left(\frac{\kappa N}{2}\right)}{\cot\left(\frac{\kappa N}{2}\right) + \cot\left(\frac{\kappa}{2}\right)} \\ &= -\frac{1}{\sin^2\left(\frac{\kappa N}{2}\right) \left(\cot\left(\frac{\kappa N}{2}\right) + \cot\left(\frac{\kappa}{2}\right) \right)} \end{aligned}$$

and therefore

$$\mathcal{E}(N+1) - \mathcal{E}(N) = \frac{4}{\pi^2 \epsilon} - \frac{\epsilon}{2(\Delta x)^2} \frac{1}{\sin^2\left(\frac{\kappa N}{2}\right) \left(\cot\left(\frac{\kappa N}{2}\right) + \cot\left(\frac{\kappa}{2}\right) \right)} \sin(\kappa). \quad (\text{A.21})$$

This expression is initially decreasing for small N and will eventually start increasing as N passes the optimal value of N whose zero, treating N as continuous, would be achieved ² for (see the discussion preceding Equation (6.40) for details on the elimination of the \tan^{-1})

$$N_{cont} = \frac{2}{\kappa} \tan^{-1} \left(\sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1} \right) = \frac{\pi}{\kappa} - 1. \quad (\text{A.22})$$

²While this can be shown more rigorously treating N as continuous and considering $\frac{\partial^2 \mathcal{E}}{\partial N^2}$, the monotonicity is intuitively obvious as very thin interfaces will spread due to the very high gradient energy density whereas very broad interfaces will contract due to the high bulk energy density. Treating N as continuous, both the actual zero of the difference and the sign of the second derivative can be found by using tangent half-angle substitution $t = \tan\left(\frac{\kappa N}{2}\right)$ and basic trigonometric identities to be given by $t = \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}$, and therefore the expression in Equation (A.22). Since the correct ‘‘guess’’ is already available, this calculation is skipped here for shortness.

In fact, using Equations (A.6) and (A.7) as well as $\kappa N_{cont} = \pi - \kappa$, one has

$$\sin^2\left(\frac{\kappa N_{cont}}{2}\right) = \frac{1 - \cos(\kappa N_{cont})}{2} = \frac{1 - \cos(\pi - \kappa)}{2} = \frac{1 + \cos(\kappa)}{2} = 1 - \frac{4(\Delta x)^2}{\pi^2 \epsilon^2} = \frac{4(\Delta x)^2}{\pi^2 \epsilon^2} \left(\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1 \right)$$

and $\cot\left(\frac{\kappa}{2}\right) = \frac{1 + \cos(\kappa)}{\sin(\kappa)} = \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}$ and by Equation (A.22) $\cot\left(\frac{\kappa N_{cont}}{2}\right) = \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}^{-1}$. It follows that the denominator in the second term in Equation (A.21) reduces to

$$\frac{4(\Delta x)^2}{\pi^2 \epsilon^2} \left(\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1 \right) \left(\sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}^{-1} + \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1} \right)$$

Combining this with Equation (A.7), one obtains

$$\begin{aligned} \mathcal{E}(N_{cont} + 1) - \mathcal{E}(N_{cont}) &= \frac{4}{\pi^2 \epsilon} - \frac{\epsilon}{2(\Delta x)^2} \frac{\frac{8(\Delta x)^2}{\pi^2 \epsilon^2} \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}}{\frac{4(\Delta x)^2}{\pi^2 \epsilon^2} \left(\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1 \right) \left(\sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1}^{-1} + \sqrt{\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1} \right)} \\ &= \frac{4}{\pi^2 \epsilon} - \frac{\frac{4}{\pi^2 \epsilon}}{\frac{4(\Delta x)^2}{\pi^2 \epsilon^2} \left(1 + \left(\frac{\pi^2 \epsilon^2}{4(\Delta x)^2} - 1 \right) \right)} = 0. \end{aligned}$$

From the global point of view, it follows that the optimal N for enforcing a transition region compatible with the first-order necessary condition is given by $\lceil N_{cont} \rceil = \lceil \frac{\pi}{\kappa} - 1 \rceil$ if $\frac{\pi}{\kappa}$ is not integer since the lower choice $\lfloor N_{cont} \rfloor$ is still in the region where the energy is smaller when adding an additional point. In contrast, if $\frac{\pi}{\kappa}$ is integer, the energies for both $N = \lfloor N_{cont} \rfloor$ and $N = \lceil N_{cont} \rceil$ coincide, meaning that from a discrete perspective (in N), the choice is indeterminate since both are energetically optimal.

A.3 The Local Analysis and the Second-Order Conditions

While the analysis in the previous section was primarily concerned with the analysis of the energetics in terms of the discrete parameter N , this does not exclude the possibility of several local minima for two successive values of N . This section is therefore more focused on the local stability in terms of continuous variations of the ϕ_i .

As recalled in Section A.1, the two values $N_{min} := \lceil \frac{\pi}{\kappa} - 1 \rceil$ and $N_{max} = \lceil \frac{\pi}{\kappa} \rceil$ together with the profile in Equation (A.5) define two admissible solutions to the first-order necessary condition for a local minimizer of the phasefield equation. While it was seen in the previous section that the lower choice $N = N_{min}$ usually is the one with the lower total energy (unless $\frac{\pi}{\kappa}$ is integer, in which case both energies are equal), this does neither imply that the second choice $N = N_{max}$ is not a local minimizer nor, without further argument, that the profile for $N = N_{min}$ is in fact also a local minimizer and not just a critical point of the energy \mathcal{E} . Investigating this question requires including second-order information. Even though the question of second-order necessary and sufficient conditions for equality- and inequality-constrained problems is a well-studied one (see e.g. [13] and [46]), due to the fact that the energy \mathcal{E} is a quadratic form in ϕ and that the bound-constraint $0 \leq \phi_i \leq 1$ is the only relevant restriction on ϕ it seems preferable to argue directly on the equation itself instead of using theorems primarily aimed at more difficult settings. The following argument is essentially a slight variation and extension of the proof of the closely related Theorem 16.4 in [54] for quadratic programming problems with linear side constraints in order to avoid the convexity assumption on the energy in that theorem³.

³As there is a large body of (primarily algorithmic) literature on such indefinite quadratic programming problems, one could certainly also find theorems specifically adapted to this situation. Since the argument below is quite instructive for the particular problem considered here, in particular for $\frac{\pi}{\kappa}$ integer, no particular effort has been made though for finding a specific reference where the required conclusions are stated in a simple and explicit form.

As the phasefield energy \mathcal{E} is quadratic in ϕ , the quadratic expansion

$$\mathcal{E}(\phi + \delta\phi) = \mathcal{E}(\phi) + \mathcal{E}'(\phi) \cdot \delta\phi + \frac{1}{2} \delta\phi \cdot \mathcal{E}''(\phi) \cdot \delta\phi \quad (\text{A.23})$$

is exact around any given value of ϕ and for any direction $\delta\phi$. $\mathcal{E}'(\phi)$ is precisely the left-hand side of by Equation (6.38), such that, if the ϕ -profile solves the first-order necessary condition, one has

$$\mathcal{E}'(\phi) \cdot \delta\phi = \sum_i \left(-2\gamma\epsilon \frac{\phi_{i+1} - 2\phi_i + \phi_{i-1}}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon} \gamma(1 - 2\phi_i) \right) \delta\phi_i = \sum_i (\mu_i^- - \mu_i^+) \delta\phi_i,$$

together with the complementarity conditions $\mu_i^- \geq 0$, $\mu_i^- = 0$ if $\phi_i > 0$ resp. $\mu_i^+ \geq 0$, $\mu_i^+ = 0$ if $\phi_i = 1$. From this, $\mathcal{E}(\phi + \delta\phi)$ in any such point can be rewritten as

$$\mathcal{E}(\phi + \delta\phi) = \mathcal{E}(\phi) + \sum_i (\mu_i^- - \mu_i^+) \delta\phi_i + \frac{1}{2} \delta\phi \cdot (\mathcal{E}''(\phi) \delta\phi), \quad (\text{A.24})$$

where $\mathcal{E}''(\phi) =: \mathbf{A}$ is the (constant) matrix characterized by the homogeneous form

$$(\mathbf{A}\psi)_i = -2\gamma\epsilon \frac{\psi_{i+1} - 2\psi_i + \psi_{i-1}}{(\Delta x)^2} - \frac{32}{\pi^2\epsilon} \gamma\psi_i. \quad (\text{A.25})$$

In order to show that a given phasefield profile satisfying the FONC is indeed a local minimizer (not necessarily strict), what needs to be shown is that

$$\mathcal{E}(\phi + t\mathbf{d}) - \mathcal{E}(\phi) = t \sum_i (\mu_i^- - \mu_i^+) d_i + \frac{1}{2} t^2 \mathbf{d} \cdot (\mathcal{E}''(\phi) \mathbf{d}) \geq 0 \quad (\text{A.26})$$

for all admissible directions \mathbf{d} and $t > 0$ sufficiently small. A first important observation is that, if any multiplier μ_i^\pm is strictly positive (i.e. the constraint is strongly active), it suffices to consider directions \mathbf{d} such that $d_i = 0$ for any such i . In fact, assuming e.g. $\mu_i^- > 0$, the complementarity condition enforces that $\phi_i = 0$ and thus the only admissible directions are such that $d_i \geq 0$. Similarly, d_i is necessarily non-positive if $\mu_i^+ > 0$. From this, it follows that $(\mu_i^- - \mu_i^+) d_i \geq 0$, and, if an entry d_i for any strongly active constraint is non-zero, will actually be strictly positive. Combined with the quadratic term being $\mathcal{O}(t^2)$, any such direction leads to a strict increase in the energy for sufficiently small $t > 0$, regardless of the values of the remaining d_j , $j \neq i$.

The only admissible directions in which the energy could potentially locally decrease are thus the ones which can vary freely if $0 < \phi_i < 1$ and, if there are weakly active constraint with $\phi_i = 0$ and $\mu_i^- = 0$ or $\phi_i = 1$ and $\mu_i^+ = 0$, the respective sign-restriction at these points. In addition, since the linear term $\sum_i (\mu_i^- - \mu_i^+) d_i$ vanishes for all such directions, one has

$$\mathcal{E}(\phi + t\mathbf{d}) - \mathcal{E}(\phi) = t^2 \mathbf{d} \cdot (\mathcal{E}''(\phi) \mathbf{d}), \quad (\text{A.27})$$

i.e. the sign of the energy difference depends only upon that of the quadratic term in \mathbf{d} .

All constraints in the bulk regions $i < 0$ and $i > N$ are strictly active (with the multiplier taking the value $\frac{16\gamma}{\pi^2\epsilon}$). Furthermore, from the derivation of the lower bound on N in Equation (6.40), the multipliers μ_0^- resp. μ_N^+ at the outermost points are strictly positive provided $N > \frac{\pi}{\kappa} - 1$ and only become zero if $N = \frac{\pi}{\kappa} - 1$. One therefore has to distinguish two cases:

1. If $\frac{\pi}{\kappa}$ is not an integer, both by N_{min} and N_{max} are strictly larger than $\frac{\pi}{\kappa} - 1$ and all active constraints are strongly active. From this, it follows that it suffices to focus on variations within the inner interface region, i.e. the interval $1 \leq i \leq N - 1$, with all relevant search directions \mathbf{d} satisfying $d_i = 0$, for $i \leq 0$ and $i \geq N$ and d_i arbitrary for $1 \leq i \leq N - 1$. In

combination with Equation (A.27), this reduces the question of being a local minimizer to the subblock of the matrix $\mathbf{A} = \mathcal{E}''(\phi)$ corresponding to values $1 \leq i \leq N - 1$ being positive (semi-)definite, resp. by Equation (A.25), to all eigenvalues of the linear difference equation

$$-2\gamma\epsilon \frac{\psi_{i+1} - 2\psi_i + \psi_{i-1}}{(\Delta x)^2} - \frac{32}{\pi^2\epsilon} \gamma \psi_i \quad , 1 \leq i \leq N - 1,$$

subject to the homogeneous Dirichlet boundary condition $\psi_0 = \psi_N = 0$ being non-negative. It is well-known that the eigenfunctions for such an operator are given by $\psi_i = c \sin\left(\frac{k\pi}{N}i\right)$, $1 \leq i \leq N$ with $1 \leq k \leq N - 1$. Using $\sin\left(\frac{k\pi}{N}(i \pm 1)\right) = \sin\left(\frac{k\pi}{N}i\right) \cos\left(\frac{k\pi}{N}\right) \pm \cos\left(\frac{k\pi}{N}i\right) \sin\left(\frac{k\pi}{N}\right)$, it is easy to see that, for each k , the corresponding eigenvalue λ_k is given by

$$\lambda_k = 4\gamma\epsilon \frac{1 - \cos\left(\frac{k\pi}{N}\right)}{(\Delta x)^2} - \frac{32}{\pi^2\epsilon} \gamma.$$

Since $0 < \frac{k\pi}{N} < \pi$ for $1 \leq k \leq N - 1$ and the cosine is decreasing on this interval, the lowest eigenvalue is given by $4\gamma\epsilon \frac{1 - \cos\left(\frac{\pi}{N}\right)}{(\Delta x)^2} - \frac{32}{\pi^2\epsilon} \gamma = \frac{4\gamma\epsilon}{(\Delta x)^2} \left(1 - \cos\left(\frac{\pi}{N}\right) - \frac{8(\Delta x)^2}{\pi^2\epsilon^2}\right)$, which is non-negative if

$$\cos\left(\frac{\pi}{N}\right) \leq 1 - \frac{8(\Delta x)^2}{\pi^2\epsilon^2} = \cos(\kappa).$$

On the relevant domain, $\cos\left(\frac{\pi}{N}\right)$ is monotonically decreasing with N , from which it follows that $\lambda_1 \geq 0$ if $\frac{\pi}{N} \geq \kappa$ i.e. if $N \leq \frac{\pi}{\kappa}$. Since $\frac{\pi}{\kappa}$ was assumed non-integer, this inequality is strictly satisfied for $N = N_{min} = \lceil \frac{\pi}{\kappa} - 1 \rceil$, whereas it is strictly violated for $N = N_{max} = \lfloor \frac{\pi}{\kappa} \rfloor$. It follows that the profile for $N = N_{min}$ is a strict local (and actually also global under the ‘‘constraint’’ of having an actual interface by the analysis in the last section) minimizer of \mathcal{E} . In contrast, the profile obtained for $N = N_{max}$ is not a local minimizer since adding an arbitrarily small multiple of the first eigenfunction to the profile - and thus in particular breaking its symmetry - will lead to a strict decrease of the energy.

2. If $\frac{\pi}{\kappa}$ is integer, the situation is somewhat more complex but also considerable more interesting. In this case, the two possible choices for N are given by $N_{min} = \lceil \frac{\pi}{\kappa} - 1 \rceil = \frac{\pi}{\kappa} - 1$ and $N_{max} = \lfloor \frac{\pi}{\kappa} \rfloor = \frac{\pi}{\kappa}$. For the latter choice, both relevant multipliers μ_0^- and μ_N^+ are strictly positive with a value given, by Equation (6.39) combined with $\cot\left(\frac{\pi}{2}\right) = 0$ and equation (A.6) by

$$\mu_0^- = \mu_N^+ = \gamma\epsilon \frac{\cos(\kappa) - 1}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon^2} \gamma = \gamma\epsilon \frac{-\frac{8(\Delta x)^2}{\pi^2\epsilon}}{(\Delta x)^2} + \frac{16}{\pi^2\epsilon^2} \gamma = \frac{8}{\pi^2\epsilon^2} \gamma,$$

i.e. a value which corresponds precisely to half the for the multipliers in the bulk. As the constraints are therefore strictly active, one can apply the same analysis as for the previous case to conclude that the profile for $N = N_{max}$ is in fact a local minimum, but with the smallest eigenvalue of the corresponding eigenvalue problem being zero due to $\cos\left(\frac{\pi}{N}\right) = \cos(\kappa)$.

In contrast, for the choice $N = \frac{\pi}{\kappa} - 1$, $\mu_0^- = \mu_N^+ = 0$. This implies that the values ϕ_0 and ϕ_N can also be varied without a first-order increase in the energy, but the entries d_0 and d_N of any admissible direction do have to satisfy a sign-restriction in order to remain compatible with the box-constraints. Here it is convenient to separately consider two types of search directions.

- The first one is formed by those directions for which at most one of the outermost interface points ϕ_0 or ϕ_N would move away from the respective constraint, i.e. directions with either $d_0 \geq 0$ and $d_N = 0$ or $d_0 = 0$ and $d_N \leq 0$. These directions are a subset of all the directions \mathbf{d} for which at most $N_{min} + 1 = N_{max} = \frac{\pi}{\kappa}$ values are free

to vary, i.e. of those directions where either $d_0 = d_{N+1} = 0$ but d_i , $1 \leq i \leq N$ can be chosen freely or where d_i , $0 \leq i \leq N-1$ can be chosen freely but $d_{-1} = d_N = 0$. This “relaxed” problem is easily seen to reduce precisely to the eigenvalue analysis performed for $N = N_{max}$ and it follows that, except for those directions generated by negative multiples of $\sin\left(\frac{\pi}{N_{max}}i\right)$ (compatible with $d_0 = 0$ and $d_{N_{min}} = d_{N_{max}-1} < 0$) or positive multiples of $\sin\left(\frac{\pi}{N_{max}}(i-1)\right)$ (compatible with $d_0 > 0$ and $d_{N_{min}+1} = d_{N_{max}} = 0$) of the eigenfunction $\sin\left(\frac{\pi}{N_{max}}i\right)$ resp. its translate lead to a strict increase in the energy. As these particular eigenfunctions leave the energy invariant, it follows that $\mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}) \geq 0$ for all such directions.

- It remains to verify that it is not possible to decrease the energy by choosing a direction with both $d_0 > 0$ and $d_{N_{min}} < 0$ corresponding to a simultaneous expansion of the inner interface in two directions. This cannot be verified using the same argument, as fixing only $d_i = 0$ for $i \leq -1$ and $i \geq N_{min} + 1 = N_{max}$ would lead to an eigenvalue problem on a domain which is “broader” than the maximal stable number of points given by N_{max} and would actually allow for an eigenfunction $\sin\left(\frac{\pi}{N_{max}+1}(i+1)\right)$, $-1 \leq i \leq N_{max}$ with a negative eigenvalue. Here the constraints on d_0 and $d_{N_{min}}$ are decisive, as these essentially exclude this type of direction since it either has to increase or decrease both ϕ_0 and $\phi_{N_{min}}$ and is therefore always incompatible with the one of the sign restriction.

More precisely, from the analysis for $N = N_{max}$, the function $\hat{d}_i = \sin\left(\frac{\pi}{N_{max}}i\right)$, $0 \leq i \leq N_{max}$ is an eigenfunction with eigenvalue 0 for the restriction of the operator $\mathcal{E}''(\phi)$ to the range of indices $1 \leq i \leq N_{max} - 1$. In contrast, when broadening the region of interest by increasing the inner interval under consideration to $0 \leq i \leq N_{max}$ and extending this function (still denoted by $\hat{\mathbf{d}}$) with $\hat{d}_{-1} = 0$, this function is not an eigenfunction anymore, but satisfies (similarly to the homogeneous form of Equation (6.39) defining μ_0^-)

$$\begin{aligned} (\mathcal{E}''(\phi)\hat{\mathbf{d}})_{i, 1 \leq i \leq N_{max}-1} &= 0 \quad \text{and} \\ (\mathcal{E}''(\phi)\hat{\mathbf{d}})_{i, 1 \leq i \leq N_{max}-1} &= 2\gamma\epsilon \frac{d_1}{(\Delta x)^2} = 2\gamma\epsilon \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2}, \end{aligned} \tag{A.28}$$

i.e. the homogeneous difference operator applied to this extension of the eigenfunction still leaves the previous (for $N = N_{min}$) interior of the interface unaffected, but leads a positive contribution to the cell $i = 0$. An arbitrary search direction $\mathbf{d} = (d_i)_{0 \leq i \leq N_{min}}$ can be decomposed as

$$d_i = \underbrace{d_i - \frac{d_{N_{min}}}{\sin\left(\frac{\pi}{N_{max}}N_{min}\right)}\hat{d}_i}_{=: d'_i} + \frac{d_{N_{min}}}{\sin\left(\frac{\pi}{N_{max}}N_{min}\right)}d_i = d'_i - \frac{d_{N_{min}}}{\sin\left(\frac{\pi}{N_{max}}\right)}\hat{d}_i,$$

where use was made of $\sin\left(\frac{\pi}{N_{max}}N_{min}\right) = \sin\left(\frac{\pi}{N_{max}}(N_{max}-1)\right) = \sin\left(\pi - \frac{\pi}{N_{max}}\right) = \sin\left(\frac{\pi}{N_{max}}\right)$ and where $d'_{N_{min}} = 0$ by the choice of prefactor for $\hat{\mathbf{d}}$.

Since $d_{N_{min}}$ is negative by assumption and $\sin\left(\frac{\pi}{N_{max}}\right)$ is positive, this decomposes any such direction \mathbf{d} as the sum of one vector \mathbf{d}' with $d'_{N_{min}} = 0$, $d'_0 = d_0 + c \sin\left(\frac{\pi}{N_{max}}\right) > d_0 > 0$ and a positive multiple $c = -\frac{d_{N_{min}}}{\sin\left(\frac{\pi}{N_{max}}\right)}$ of the extension of the eigenvector $\hat{\mathbf{d}}$.

From Equation (A.28), it further follows that

$$\begin{aligned} (\mathcal{E}''(\phi)\mathbf{d})_{1 \leq i \leq N_{min}} &= (\mathcal{E}''(\phi)\mathbf{d}')_{1 \leq i \leq N_{min}} \quad \text{and} \\ (\mathcal{E}''(\phi)\mathbf{d})_0 &= (\mathcal{E}''(\phi)\mathbf{d}')_0 + 2\gamma\epsilon c \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2} \end{aligned} \quad (\text{A.29})$$

since $\hat{\mathbf{d}}$ does not affect the results of the operator for any index $i > 0$. Taking the scalar product with \mathbf{d} , one has

$$\begin{aligned} \mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}) &= d_0 (\mathcal{E}''(\phi)\mathbf{d})_0 + \sum_{i=1}^{N_{min}} d_i (\mathcal{E}''(\phi)\mathbf{d})_i \\ &= d_0 \left((\mathcal{E}''(\phi)\mathbf{d}')_0 + 2\gamma\epsilon c \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2} \right) + \sum_{i=1}^{N_{min}} d_i (\mathcal{E}''(\phi)\mathbf{d}')_i \\ &= \underbrace{2\gamma\epsilon c \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2} d_0}_{>0} + \mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}') > \mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}'). \end{aligned}$$

Again with Equation (A.29), $\mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}')$ can be expanded in a similar fashion using the symmetry of $\mathcal{E}''(\phi)$ as

$$\begin{aligned} \mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}') &= (\mathcal{E}''(\phi)\mathbf{d}) \cdot \mathbf{d}' = (\mathcal{E}''(\phi)\mathbf{d})_0 d'_0 + \sum_{i=1}^{N_{min}} (\mathcal{E}''(\phi)\mathbf{d})_i d'_i \\ &= \left((\mathcal{E}''(\phi)\mathbf{d}')_0 + 2\gamma\epsilon c \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2} \right) d'_0 + \sum_{i=1}^{N_{min}} (\mathcal{E}''(\phi)\mathbf{d}')_i d'_i. \end{aligned}$$

The term $2\gamma\epsilon c \frac{\sin\left(\frac{\pi}{N_{max}}\right)}{(\Delta x)^2} d'_0$ is again strictly positive, such that $\mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d})$ can finally be estimated from below as $\mathbf{d} \cdot (\mathcal{E}''(\phi)\mathbf{d}) > \mathbf{d}' \cdot (\mathcal{E}''(\phi)\mathbf{d}')$. Since the last entry of \mathbf{d}' is in addition zero by construction, the sign of this term in fact only depends on the sign of this quadratic form restricted to the subvector $(\mathbf{d}')_{0 \leq i \leq N_{min}-1}$, i.e. on vectors which have at most one additional non-zero entry as compared to the original inner transition region. This is precisely the question already investigated in the previous point, finally showing that the “standard” profile with $N = N_{min}$ is also a local (non-strict) minimum.

Bibliography

- [1] Patrick Altschuh. Skalenübergreifende analyse makroporöser membranen im kontext digitaler zwillinge, 2020. Dissertation, Karlsruher Institut für Technologie (KIT).
- [2] L. Ambrosio, N. Fusco, and D. Pallara. *Functions of Bounded Variation and Free Discontinuity Problems*. Oxford Mathematical Monographs. Clarendon Press, Oxford, 2000.
- [3] Kais Ammar, Benoît Appolaire, Georges Cailletaud, and Samuel Forest. Combining phase field approach and homogenization methods for modelling phase transformation in elastoplastic media. *European Journal of Computational Mechanics/Revue Européenne de Mécanique Numérique*, 18(5-6):485–523, 2009.
- [4] PG Kubendran Amos, Ephraim Schoof, Daniel Schneider, and Britta Nestler. Chemo-elastic phase-field simulation of the cooperative growth of mutually-accommodating widmanstätten plates. *Journal of Alloys and Compounds*, 767:1141–1154, 2018.
- [5] PG Kubendran Amos, Ephraim Schoof, Nick Streichan, Daniel Schneider, and Britta Nestler. Phase-field analysis of quenching and partitioning in a polycrystalline fe-c system under constrained-carbon equilibrium condition. *Computational Materials Science*, 159:281–296, 2019.
- [6] Hedy Attouch, Giuseppe Buttazzo, and Gérard Michaille. *Variational analysis in Sobolev and BV spaces: applications to PDEs and optimization*. SIAM, Philadelphia, Pa., 2. ed. edition, 2014.
- [7] Jean-Pierre Aubin and Hélène Frankowska. *Set-valued analysis*. Birkhäuser, Boston, Mass., 2009.
- [8] H. Bauschke and P. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. CMS Books in Mathematics. Springer, Cham, 2 edition, 2017.
- [9] Marouen Ben Said, Michael Selzer, Britta Nestler, Daniel Braun, Christian Greiner, and Harald Garcke. A phase-field approach for wetting phenomena of multiphase droplets on solid surfaces. *Langmuir*, 30(14):4033–4039, 2014.
- [10] Benzi, Michele and Golub, Gene H and Liesen, Jörg. Numerical solution of saddle point problems. *Acta numerica*, 14:1–137, 2005.
- [11] Bertsekas, Dimitri P. *Nonlinear programming*. Athena scientific, Belmont, Mass., 1999.
- [12] L. Blank et al. Allen-Cahn and Cahn-Hilliard variational inequalities solved with optimization techniques. *International Series of Numerical Mathematics*, 160:21–35, 2012.
- [13] F. B. Bonnans et al. *Numerical Optimization - Theoretical and Practical Aspects*. Universitext. Springer, Berlin, 2 edition, 2006.

- [14] J. F. Bonnans and A. Shapiro. *Perturbation analysis of optimization problems*. Springer Series in Operations Research. Springer, New York, 2000.
- [15] A. Braides. *Gamma-convergence for beginners*, volume 22 of *Oxford lecture series in mathematics and its applications*. Oxford University Press, Oxford, 2002.
- [16] Haim Brezis. *Functional analysis, Sobolev spaces and partial differential equations*. Springer, New York, NY, 2010.
- [17] Céa, Jean and Glowinski, Roland. Sur des méthodes d’optimisation par relaxation. *Revue française d’automatique informatique recherche opérationnelle. Mathématique*, 7(R3):5–31, 1973.
- [18] Armen G. Chačaturijan. *Theory of structural transformations in solids*. Wiley, New York [u.a.], 1983.
- [19] Abhik Choudhury and Britta Nestler. Grand-potential formulation for multicomponent phase transformations combined with thin-interface asymptotics of the double-obstacle potential. *Physical Review E*, 85(2):021602, 2012.
- [20] Clarke, F. *Functional analysis, calculus of variations and optimal control*, volume 264. Springer Science & Business Media, London, 2013.
- [21] Klaus Deckelnick, Gerhard Dziuk, and Charles M Elliott. Computation of geometric partial differential equations and mean curvature flow. *Acta numerica*, 14:139, 2005.
- [22] N Dinh and V Jeyakumar. Farkas’ lemma: three decades of generalizations for mathematical optimization. *Top*, 22(1):1–22, 2014.
- [23] A Durga, Patrick Wollants, and Nele Moelans. Evaluation of interfacial excess contributions in different phase-field models for elastically inhomogeneous systems. *Modelling and simulation in materials science and engineering*, 21(5):055018, 2013.
- [24] A Durga, Patrick Wollants, and Nele Moelans. A quantitative phase-field model for two-phase elastically inhomogeneous systems. *Computational Materials Science*, 99:81–95, 2015.
- [25] J Eiken, B Böttger, and I Steinbach. Multiphase-field approach for multicomponent alloys with extrapolation scheme for numerical application. *Physical review E*, 73(6):066122, 2006.
- [26] I. Ekeland and R. Témam. *Convex analysis and variational problems*. Classics in Applied Mathematics. SIAM, Philadelphia, Pa., 1999.
- [27] F. Facchinei and J.S. Pang. *Finite-Dimensional Variational Inequalities and Complementary Problems*, volume 1 of *Springer Series in Operations Research*. Springer, New York, 2007.
- [28] Avner Friedman. *Variational Principles and Free Boundary Problems*. Dover Publications, Mineola, NY, dover ed. edition, 2010.
- [29] H. Garcke et al. Allen-Cahn systems with volume constraints. *Mathematical Models and Methods in Applied Sciences*, 18(8):1347–1381, 2008.
- [30] C. Gräser. Convex minimization and phase field models, 2011. Dissertation Freie Universität Berlin.
- [31] C. Gräser, R. Kornhuber, and U. Sack. Time discretizations of anisotropic Allen-Cahn equations. *IMA Journal of Numerical Analysis*, 33(48):1226–1244, 2013.

- [32] Carsten Gräser and Ralf Kornhuber. On preconditioned Uzawa-type iterations for a saddle point problem with inequality constraints. In *Domain decomposition methods in science and engineering XVI*, pages 91–102. Springer, Berlin, Heidelberg, 2007.
- [33] Carsten Gräser and Oliver Sander. Truncated nonsmooth Newton multigrid methods for simplex-constrained minimization problems. *Preprint 384, IGPM Aachen*, 2014.
- [34] P. Grisvard. *Elliptic Problems in Nonsmooth Domains*, volume 69 of *Classics in Applied Mathematics*. SIAM, Philadelphia, Pa., 2011.
- [35] Christoph Herrmann, Ephraim Schoof, Daniel Schneider, Felix Schwab, Andreas Reiter, Michael Selzer, and Britta Nestler. Multiphase-field model of small strain elasto-plasticity according to the mechanical jump conditions. *Computational Mechanics*, 62(6):1399–1412, 2018.
- [36] Johannes Hötzer. *Massiv-parallele und großskalige Phasenfeldsimulationen zur Untersuchung der Mikrostrukturentwicklung*. Schriftenreihe des Instituts für Angewandte Materialien ; Band 70. KIT Scientific Publishing, Karlsruhe, 2017. Dissertation, Karlsruher Institut für Technologie (KIT).
- [37] K. Ito and K. Kunisch. *Lagrange Multiplier Approach to Variational Problems and Applications*. Advances in Design and Control. SIAM, Philadelphia, Pa., 2008.
- [38] Boško S Jovanović and Endre Süli. *Analysis of finite difference schemes: for linear partial differential equations with generalized solutions*, volume 46. Springer Science & Business Media, London, 2013.
- [39] Michael Kellner. Modellierung mehrkomponentiger Materialsysteme für die Phasenfeldmethode und Analyse der simulierten Mikrostrukturen, 2020. Dissertation, Karlsruher Institut für Technologie (KIT).
- [40] Seong Gyoon Kim, Dong Ik Kim, Won Tae Kim, and Yong Bum Park. Computer simulations of two-dimensional and three-dimensional ideal grain growth. *Physical Review E*, 74(6):061605, 2006.
- [41] Kim, Seong Gyoon and Kim, Won Tae and Suzuki, Toshio. Interfacial compositions of solid and liquid in a phase-field model with finite interface thickness for isothermal solidification in binary alloys. *Physical Review E*, 58(3):3316–3323, 1998.
- [42] Kim, Seong Gyoon and Kim, Won Tae and Suzuki, Toshio. Phase-field model for binary alloys. *Physical Review E*, 60(6):7186–7197, 1999.
- [43] D. Kinderlehrer and G. Stampacchia. *An Introduction to Variational Inequalities and Their Applications*. Classics in Applied Mathematics. SIAM, Philadelphia, Pa., 2000.
- [44] Cornelius Lanczos. *The variational principles of mechanics*. Dover, New York, 2012.
- [45] Stephan Luckhaus and Luciano Modica. The gibbs-thompson relation within the gradient theory of phase transitions. *Archive for Rational Mechanics and Analysis*, 107(1):71–83, 1989.
- [46] Luenberger, David G and Ye, Yinyu and others. *Linear and nonlinear programming*, volume 2. Springer, Cham, 1984.
- [47] G. Dal Maso. *An Introduction to Gamma-Convergence*. Progress in Nonlinear Differential Equations and Their Applications. Birkhäuser, Boston, 1993.

- [48] W. McLean. *Strongly elliptic systems and boundary integral equations*. Cambridge University Press, Cambridge, 2000.
- [49] L. Modica. Gradient theory of phase transitions with boundary contact energy. *Annales de l'I.H.P.*, 4(5):487–512, 1987.
- [50] L. Modica. The gradient theory of phase transitions and the minimal interface criterion. *Archive for Rational Mechanics and Analysis*, 98(2):123–142, 1987.
- [51] J Mosler, O Shchyglo, and H Montazer Hojjat. A novel homogenization method for phase field approaches based on partial rank-one relaxation. *Journal of the Mechanics and Physics of Solids*, 68:251–266, 2014.
- [52] Britta Nestler, Harald Garcke, and Björn Stinner. Multicomponent alloy solidification: phase-field modeling and simulations. *Physical Review E*, 71(4):041609, 2005.
- [53] Britta Nestler, Frank Wendler, Michael Selzer, Björn Stinner, and Harald Garcke. Phase-field model for multiphase systems with preserved volume fractions. *Physical Review E*, 78(1):011604, 2008.
- [54] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research. Springer, New York, 2 edition, 2006.
- [55] James M Ortega and Werner C Rheinboldt. *Iterative solution of nonlinear equations in several variables*, volume 30. Siam, Philadelphia, Pa., 1970.
- [56] M. Plapp. Unified derivation of phase-field models for alloy solidification from a grand-potential functional. *Physical Review E*, 84(031601):1–15, 2011.
- [57] R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media, Berlin, 2009.
- [58] Ruszczyński, Andrzej P and Ruszczyński, Andrzej. *Nonlinear optimization*, volume 13. Princeton university press, Princeton, 2006.
- [59] Alexander A Samarskii. *The theory of difference schemes*, volume 240. CRC Press, New York, 2001.
- [60] Lavinia Sarbu. *Primal-dual active set methods for Allen-Cahn variational inequalities*. PhD thesis, University of Sussex, 2010.
- [61] Daniel Schneider. Phasenfildmodellierung mechanisch getriebener grenzflächenbewegungen in mehrphasigen systemen, 2016. Dissertation, Karlsruher Institut für Technologie (KIT).
- [62] Daniel Schneider, Ephraim Schoof, Oleg Tschukin, Andreas Reiter, Christoph Herrmann, Felix Schwab, Michael Selzer, and Britta Nestler. Small strain multiphase-field model accounting for configurational forces and mechanical jump conditions. *Computational Mechanics*, 61(3):277–295, 2018.
- [63] Daniel Schneider, Felix Schwab, Ephraim Schoof, Andreas Reiter, Christoph Herrmann, Michael Selzer, Thomas Böhlke, and Britta Nestler. On the stress calculation within phase-field approaches: a model for finite deformations. *Computational Mechanics*, 60(2):203–217, 2017.
- [64] Daniel Schneider, Oleg Tschukin, Abhik Choudhury, Michael Selzer, Thomas Böhlke, and Britta Nestler. Phase-field elasticity model based on mechanical jump conditions. *Computational Mechanics*, 55(5):887–901, 2015.

- [65] Ephraim Schoof. Chemomechanische modellierung der wärmebehandlung von stählen mit der phasenfeldmethode, 2020. Dissertation, Karlsruher Institut für Technologie (KIT).
- [66] Ephraim Schoof, Christoph Herrmann, Nick Streichhan, Michael Selzer, Daniel Schneider, and Britta Nestler. On the multiphase-field modeling of martensitic phase transformation in dual-phase steel using j2-viscoplasticity. *Modelling and Simulation in Materials Science and Engineering*, 27(2):025010, 2019.
- [67] Michael Selzer. Mechanische und strömungsmechanische topologieoptimierung mit der phasenfeldmethode, 2014. Dissertation, Karlsruher Institut für Technologie (KIT).
- [68] I Steinbach and F Pezzolla. A generalized field method for multiphase transformations using interface fields. *Physica D: Nonlinear Phenomena*, 134(4):385–393, 1999.
- [69] Ingo Steinbach and Markus Apel. Multi phase field model for solid state transformation with elastic strain. *Physica D: Nonlinear Phenomena*, 217(2):153–160, 2006.
- [70] Philipp Steinmetz. Simulation der bei der gerichteten erstarrung ternärer eutektika entstehenden mikrostruktur mit der phasenfeldmethode, 2017. Dissertation, Karlsruher Institut für Technologie (KIT).
- [71] Bjorn Stinner, Britta Nestler, and Harald Garcke. A diffuse interface model for alloys with multiple components and phases. *SIAM Journal on Applied Mathematics*, 64(3):775–799, 2004.
- [72] Ying Sun and Christoph Beckermann. Sharp interface tracking using the phase-field equation. *Journal of Computational Physics*, 220(2):626–653, 2007.
- [73] R. Temam. *Navier-Stokes Equations - Theory and Numerical Analysis*. AMS Chelsea Publishing, Providence, 2001.
- [74] Oleg Tschukin. Phase-field modelling of welding and of elasticity-dependent phase transformations, 2017. Dissertation, Karlsruher Institut für Technologie (KIT).
- [75] Oleg Tschukin, Alexander Silberzahn, Michael Selzer, Prince GK Amos, Daniel Schneider, and Britta Nestler. Concepts of modeling surface energy anisotropy in phase-field approaches. *Geothermal Energy*, 5(1):19, 2017.
- [76] R. S. Varga. *Matrix Iterative Analysis*, volume 27 of *Springer Series in Computational Mathematics*. Springer, Berlin, 2000.
- [77] P. S. Vassilevski. *Multilevel Block Factorization Preconditioners*. Springer, New York, 2008.
- [78] Alexander Vondrous. Grain growth behavior and efficient large scale simulations of recrystallization with the phase-field method, 2014. Dissertation, Karlsruher Institut für Technologie (KIT).
- [79] K. Yosida. *Functional Analysis*. Classics in Mathematics. Springer, Berlin, 6 edition, 1980.
- [80] Zeidler, E. *Nonlinear functional analysis and its applications: III: variational methods and optimization*. Springer Science & Business Media, New York, 2013.