



**HAL**  
open science

## **Reactive flash for ideal multiphase mixtures: Unified formulation and efficient computation**

Maxime Jonval, Ibtihel Ben Gharbia, Clément Cancès, Thibault Faney, Quang Huy Tran

► **To cite this version:**

Maxime Jonval, Ibtihel Ben Gharbia, Clément Cancès, Thibault Faney, Quang Huy Tran. Reactive flash for ideal multiphase mixtures: Unified formulation and efficient computation. 2025. <hal-05393808>

**HAL Id: hal-05393808**

**<https://hal.science/hal-05393808v1>**

Preprint submitted on 2 Dec 2025

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

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



HAL Authorization

# Reactive flash for ideal multiphase mixtures: Unified formulation and efficient computation

Maxime Jonval\*    Ibtihel Ben Gharbia†    Clément Cancès‡    Thibault Faney†  
Quang-Huy Tran†

## Abstract

Multiphase chemical equilibrium problems lead to nonlinear systems with complementarity constraints, which become particularly challenging when phases may vanish. We introduce a new algebraic formulation of the equilibrium problem based on *extended mole fractions*, derived from the subdifferential of the Gibbs free energy, and establish its equivalence with the classical minimization problem. Our analysis provides new conditions ensuring the uniqueness of solutions, even when some phases disappear. Building on this formulation, we propose two parametrized Newton-based strategies: one reformulates the relation between species quantities and chemical potentials, while the other parametrizes the complementarity conditions directly. Numerical experiments on a system with 72 species and 22 phases confirm the robustness and efficiency of the proposed methods. In tests with randomized inputs, both strategies achieve success rates above 90% with moderate iteration counts, outperforming established approaches such as the Newton–min and Fischer–Burmeister complementarity functions, and interior-point methods.

## 1 Introduction

Simulating reactive transport remains a central challenge in many areas of science and engineering, including porous media flow, combustion, and chemical reactor design. In porous media in particular, accurate modeling of reactive transport is essential for applications such as CO<sub>2</sub> and H<sub>2</sub> storage and geothermal energy production. A key computational bottleneck in these simulations lies in solving nonlinear chemical equilibrium equations, which must be evaluated at each time step and in every mesh cell. Even modest gains in the robustness and efficiency of equilibrium solvers can therefore yield significant improvements in overall simulation performance.

Newton’s method is the standard approach for solving chemical equilibrium equations, but it faces significant challenges in multiphase settings. In earlier work [18], we investigated the single-phase case and showed that the parametrization technique proposed by Brenner and Cancès [6] can enhance robustness. In the present article, we extend this approach to the more complex case of multiphase systems.

Unlike single-phase systems, multiphase equilibrium problems are complicated by the nonnegativity constraints on species quantities, which allow certain phases to vanish at equilibrium. As a result, the set of present phases is not known a priori. Two main families of methods have been developed to address this issue. Combinatorial approaches compute equilibria for all possible subsets of phases [21, 9], but become computationally prohibitive as the number of candidate phases grows, due to the exponential increase in subsets to evaluate. Unified formulations, by contrast, introduce extended mole fractions and complementarity constraints to account for all phases simultaneously [19, 4]. By treating all phases within a single framework, unified formulations avoid combinatorial explosion and provide a more scalable approach for systems with many potential phases.

The concept of extended mole fractions was first introduced by Lauser *et al.* [19] in the context of multiphase flow in porous media, where complementarity constraints were used to handle phase appearance and disappearance within a unified system of equations. This formulation was later applied to non-reactive phase equilibrium problems by Vu *et al.* [4], who provided a detailed mathematical analysis.

---

\*Inria, Univ. Côte d’Azur, CNRS, UMR 7351 – Laboratoire J.A. Dieudonné, 06108 Nice, France. [maxime.jonval@inria.fr](mailto:maxime.jonval@inria.fr)

†IFP Energies nouvelles, 1 et 4 avenue de Bois Préau, 92852 Rueil-Malmaison Cedex, France. [ibtihel.ben-gharbia@ifpen.fr](mailto:ibtihel.ben-gharbia@ifpen.fr), [thibault.faney@ifpen.fr](mailto:thibault.faney@ifpen.fr), [quang-huy.tran@ifpen.fr](mailto:quang-huy.tran@ifpen.fr)

‡Univ. Lille, CNRS, Inria, UMR 8524 – Laboratoire Paul Painlevé, 59000 Lille, France. [clement.cances@inria.fr](mailto:clement.cances@inria.fr)

Coatléven and Michel [8] extended this approach to multiphase chemical equilibria by embedding the optimality criterion of Smith *et al.* [27] into the Karush–Kuhn–Tucker (KKT) conditions.

An alternative strategy is the species-quantity formulation proposed by Leal *et al.* [20], in which species quantities serve as the primary variables and a complementarity equation is written for each species. However, as observed by Shapiro and Shapley [26], if a species is present in a given phase, then the entire phase must also be present. This observation supports the formulation of complementarity conditions at the phase level rather than at the species level. More broadly, complementarity problems have been extensively studied, both theoretically and numerically, and provide a robust framework for modeling phase transitions and their interactions [5, 3, 2, 11].

Our contributions are threefold. First, we establish a new condition guaranteeing uniqueness even when phases vanish. Second, we derive a unified Gibbs-based formulation of multiphase chemical equilibrium and prove its equivalence with the classical minimization problem. Third, we develop two parametrized Newton strategies, including a novel complementarity parametrization method, which achieves both robustness and efficiency improvements. Numerical experiments confirm the effectiveness of these strategies across systems of up to 72 species and 22 phases.

Section 2 introduces the chemical system, defines the mass conservation laws and the Gibbs free energy, and formulates the equilibrium problem as a constrained minimization problem.

In Section 3, we extend classical results on existence and uniqueness of minimizers from Shapiro and Shapley [26], who showed that uniqueness is guaranteed when all phases are present, but not necessarily when some phases vanish. We generalize the notion of azeotropy to propose a new condition that ensures uniqueness even in the presence of phase disappearance, thereby strengthening the analytical foundation of the model. We then derive a unified formulation based on extended mole fractions and establish the uniqueness of its solution under reasonable assumptions. This formulation involves one complementarity problem per phase, which requires an adapted approach for its numerical resolution.

Section 4 is devoted to the numerical resolution of the unified formulation using Newton’s method. We employ the parametrization techniques introduced in [18] to reformulate the relation between extended mole fractions and chemical potentials. In addition, we propose a novel method for addressing the complementarity problem, the complementarity parametrization technique, which parametrizes the relationship between the complementary variables.

Finally, Section 5 presents numerical experiments validating our proposed methods. We begin with two simple test cases, followed by a more demanding benchmark involving 22 phases. This benchmark is designed to evaluate the robustness and performance of the algorithm under complex conditions. The results demonstrate the effectiveness of our numerical strategy, in particular the complementarity parametrization approach, in resolving phase presence conditions both accurately and efficiently.

## 2 Mathematical formulation for multiphase chemical equilibria

A chemical species is a group of atoms characterized by its molecular formula and the phase to which it belongs. Typical phases include aqueous (aq) solutions such as diluted electrolytes or oils, gaseous (g) phases, and solid (s) phases such as pure minerals. For instance  $\text{H}_2\text{O}(\text{aq})$ ,  $\text{H}^+(\text{aq})$ ,  $\text{NaCl}(\text{aq})$ ,  $\text{CO}_2(\text{g})$ ,  $\text{H}_2\text{O}(\text{g})$ ,  $\text{CaCO}_3(\text{s})$  are all distinct species. Note that  $\text{H}_2\text{O}(\text{aq})$  and  $\text{H}_2\text{O}(\text{g})$  are considered different species, since they belong to different phases. Similarly,  $\text{SiO}_2(\text{s})$  can occur as quartz or amorphous silica, which correspond to two distinct solid phases, denoted  $\text{SiO}_2(\text{quartz})$  and  $\text{SiO}_2(\text{amorphous} - \text{silica})$ .

For a given temperature  $T$  and pressure  $P$ , a chemical system is defined as

$$\mathcal{S}_{P,T} = \{\mathcal{P}, \mathcal{E}, \mathcal{C}, \mathcal{R}\}$$

where

- $\mathcal{P} = \{P_1, \dots, P_{N_{\text{ph}}}\}$  is the set of  $N_{\text{ph}}$  phases;
- $\mathcal{E} = \{E_1, \dots, E_M\}$  is the set of  $M$  chemical elements (typically atoms);
- $\mathcal{C} = \{C_1, \dots, C_N\}$  is the set of  $N$  chemical species, with  $N > M$ ;
- $\mathcal{R} = \{R_1, \dots, R_{N-M}\}$  is the set of  $N - M$  chemical reactions among these species.

The set  $\mathcal{E}$  contains all the elements that compose the species of the set  $\mathcal{C}$ , while the reactions in  $\mathcal{R}$  describe how these species interact. A chemical reaction  $R_j \in \mathcal{R}$  can be written as

$$\sum_{i=1}^N s_{ij} C_i = 0,$$

where the  $s_{ij}$  are the stoichiometric coefficients indicating the number of molecules of  $C_i$  involved in reaction  $R_j$ .

The set  $\mathcal{P}$  contains the phases of the system. To link each species  $C_i$  to its phase  $\alpha$ , we introduce the mapping

$$\sigma : i \in \{1, \dots, N\} \mapsto \alpha \in \{1, \dots, N_{Ph}\} \quad (1)$$

so that  $\sigma^{-1}(\alpha)$  is the set of all species in phase  $\alpha$ . We denote by  $N^\alpha = \#\sigma^{-1}(\alpha)$  the number of species in phase  $\alpha$ .

## 2.1 Mass conservation

The systems studied in this article are closed, meaning that no exchange of matter occurs with the surroundings. In such systems, the principle of mass conservation ensures that the total amount  $\mathbf{b} = (b_1, \dots, b_M)$  of each element of  $\mathcal{E}$  remains constant. Each species  $C_i$  in the set of species  $\mathcal{C}$  has an associated *formula vector*  $\mathbf{a}_i$  in the basis of elements  $\mathcal{E}$ . For example, if  $\mathcal{E} = (\text{H}, \text{C}, \text{O})$  and  $C_i = \text{HCO}_3^-$ , then  $\mathbf{a}_i = (1, 1, 3)^T$ . Based on these vectors, the set of species  $\mathcal{C}$  is partitioned into two subsets:

- Primary species  $\mathcal{C}_{Pr} = \{C_1, \dots, C_M\}$ , whose formula vectors are linearly independent. They form the primary basis of the system. The number of primary species is equal to the number of elements  $M$ .
- Secondary species  $\mathcal{C}_{Sd} = \{C_{M+1}, \dots, C_N\}$ , whose formula vectors are linear combinations of the primary species. Their number is  $N - M$ , corresponding to the  $N - M$  reactions in the set  $\mathcal{R}$ .

The choice of the primary species is not unique. For clarity, we order the species so that primary species come first, followed by secondary species. This leads to the *formula matrix*  $\mathbf{A} = [\mathbf{A}_{Pr}, \mathbf{A}_{Sd}]$  where  $\mathbf{A}_{Pr}$  is an  $M \times M$  invertible matrix formed from the primary species formula vectors, and  $\mathbf{A}_{Sd}$  is an  $M \times (N - M)$  matrix formed from the secondary species formula vectors.

**Example.** Consider a system consisting of one aqueous phase and one pure quartz mineral:

$$\begin{aligned} \mathcal{P} &= \{\text{aqueous, quartz}\} \\ \mathcal{E} &= \{\text{H, O, Si}\} \\ \mathcal{C} &= \{\text{H}_2\text{O}, \text{H}^+, \text{OH}^-, \text{SiO}_2(\text{aq}), \text{SiO}_2(\text{quartz})\} \\ \mathcal{R} &= \{\text{H}_2\text{O} = \text{H}^+ + \text{OH}^-, \text{SiO}_2(\text{aq}) = \text{SiO}_2(\text{quartz})\} \end{aligned}$$

If we choose the primary species  $\mathcal{C}_{Pr} = \{\text{H}^+, \text{OH}^-, \text{SiO}_2(\text{aq})\}$  and the secondary species  $\mathcal{C}_{Sd} = \{\text{H}_2\text{O}, \text{SiO}_2(\text{quartz})\}$ , the corresponding formula matrix is

$$\mathbf{A} = \begin{array}{ccccc} & \text{H}^+ & \text{OH}^- & \text{SiO}_2(\text{aq}) & \text{H}_2\text{O} & \text{SiO}_2(\text{quartz}) \\ \left[ \begin{array}{ccccc} 1 & 1 & 0 & 2 & 0 \\ 0 & 1 & 2 & 1 & 2 \\ 0 & 0 & 1 & 0 & 1 \end{array} \right] & \text{H} \\ & & & & & \text{O} \\ & & & & & \text{Si} \end{array}.$$

Let  $\mathbf{n} = (n_1, \dots, n_N)$  denote the molar quantities of the species in  $\mathcal{C}$ . The conservation of elements is expressed as

$$\mathbf{A}\mathbf{n} = \mathbf{b}. \quad (2)$$

Since we study multiphase equilibria, it is convenient to partition  $\mathbf{n}$  and  $\mathbf{A}$  into sub-vectors  $\mathbf{n}^\alpha$  and sub-matrices  $\mathbf{A}^\alpha$ , each corresponding to the species belonging to a given phase  $\alpha$ . In the example above,  $\mathbf{n}^{\text{aq}} = (n_{\text{H}^+}, n_{\text{OH}^-}, n_{\text{SiO}_2(\text{aq})}, n_{\text{H}_2\text{O}})$ ,  $\mathbf{n}^{\text{quartz}} = (n_{\text{SiO}_2(\text{quartz})})$ ,

$$\mathbf{A}^{\text{aq}} = \begin{bmatrix} 1 & 1 & 0 & 2 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{A}^{\text{quartz}} = \begin{bmatrix} 0 \\ 2 \\ 1 \end{bmatrix}.$$

Equation (2) then becomes

$$\sum_{\alpha=1}^{N_{Ph}} \mathbf{A}^{\alpha} \mathbf{n}^{\alpha} = \mathbf{b}.$$

We define the *stoichiometry matrix*

$$\mathbf{S} = \begin{bmatrix} \mathbf{A}_{Pr}^{-1} \mathbf{A}_{Sd} \\ -\mathbf{I}_{Sd} \end{bmatrix}. \quad (3)$$

By construction,  $\mathbf{A}\mathbf{S} = \mathbf{0}$  and  $\text{rank } \mathbf{S} = N - M$ . The entries of  $\mathbf{S}$  coincide with the stoichiometric coefficients of the reactions in  $\mathcal{R}$ . For the example system,

$$\mathbf{S} = \begin{array}{cc} \text{R}_1 & \text{R}_2 \\ \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{bmatrix} & \begin{array}{l} \text{H}^+ \\ \text{OH}^- \\ \text{SiO}_2(\text{aq}) \\ \text{H}_2\text{O} \\ \text{SiO}_2(\text{quartz}) \end{array} \end{array}.$$

The fundamental link between  $\mathbf{A}$  and  $\mathbf{S}$  is

$$\ker \mathbf{S}^T = (\ker \mathbf{A})^{\perp} = \text{Im } \mathbf{A}^T. \quad (4)$$

Moreover, the components of an element in  $\ker \mathbf{A} \setminus \{\mathbf{0}\}$  do not all have the same sign. In particular

$$\ker \mathbf{A} \cap \mathbb{R}_+^N = \{\mathbf{0}\}. \quad (5)$$

This means that no strictly positive combination of species can vanish under the conservation laws. Proofs of these properties can be found in [18], and further details on relationship between  $\mathbf{A}$  and  $\mathbf{S}$  can be found in the monograph of Smith and Missen [28].

## 2.2 Gibbs free energy and chemical potentials

The state of a closed system  $\mathcal{S}_{P,T}$  at fixed pressure  $P$  and temperature  $T$  is characterized by the Gibbs free energy function  $G : \mathbb{R}_+^N \rightarrow \mathbb{R}$ .

The Gibbs free energy function decomposes into phase-specific functions  $G_{\alpha} : \mathbb{R}_+^{N^{\alpha}} \rightarrow \mathbb{R}$  that depend solely on the species in the phase  $\alpha$ . Each  $G_{\alpha}$  is extensive, i.e. homogeneous of degree one in the molar quantities, and therefore by Euler's theorem on homogeneous functions admits the representation:

$$G(\mathbf{n}) = \sum_{\alpha=1}^{N_{Ph}} G_{\alpha}(\mathbf{n}^{\alpha}) \quad \text{where} \quad G_{\alpha}(\mathbf{n}^{\alpha}) = \sum_{i \in \sigma^{-1}(\alpha)} n_i \frac{\partial G_{\alpha}(\mathbf{n}^{\alpha})}{\partial n_i}. \quad (6)$$

The chemical potential of the species  $C_i$  is defined as  $\mu_i(\mathbf{n}^{\alpha}) = \partial G_{\alpha}(\mathbf{n}^{\alpha}) / \partial n_i = \partial G(\mathbf{n}) / \partial n_i$ . In practice, chemical potentials are commonly expressed as

$$\mu_i = \mu_i^{\circ}(P, T) + RT \ln a_i(\mathbf{n}^{\alpha}). \quad (7)$$

where  $\mu_i^{\circ}(P, T)$  is the chemical potential of species  $C_i$  in phase  $\alpha$  in its standard state at pressure  $P$  and temperature  $T$ , to be computed from thermodynamic tables, and  $a_i$  is the activity of species  $C_i$  that depends in general on the concentration of all the species in phase  $\alpha$ .

The activity of a species  $C_i$  is written as  $a_i = \gamma_i x_i$ , where  $\gamma_i$  is the activity coefficient, accounting for non-ideal interactions, and  $x_i$  denotes the mole fraction of  $C_i$  in phase  $\alpha$ , defined by

$$x_i := n_i / \sum_{j \in \sigma^{-1}(\alpha)} n_j = n_i / \langle \mathbf{n}^{\alpha}, \mathbf{1} \rangle, \quad \text{where} \quad \mathbf{1} := (1, \dots, 1)^T.$$

Numerous, increasingly complex activity models for  $\gamma_i$  are available in the literature (see *e.g.* [22, 31]). The simplest is the ideal activity model in which  $\gamma_i = 1$ . This model corresponds to an ideal mixture in which all intermolecular interactions are equivalent. In this case, the activity reduces to the mole fraction. The resulting ideal Gibbs energy in (6) is convex on  $\mathbb{R}_+^N$  (see [26, Theorem 8.13]). We restrict our attention to this case in the remainder of the article.

### 2.3 Multiphase chemical equilibrium as a minimization problem

At fixed temperature  $T$ , pressure  $P$ , and element quantities  $\mathbf{b}$ , the equilibrium state of the system is characterized by the minimization of the Gibbs free energy  $G$ . In other words, the vector of species quantities  $\mathbf{n}$  is obtained by solving a constrained optimization problem:

$$\mathbf{n} \in \arg \min \{G(\mathbf{n}) \mid \mathbf{A}\mathbf{n} = \mathbf{b} \text{ and } \mathbf{n} \geq \mathbf{0}\}. \quad (8)$$

This formulation reflects two fundamental physical principles: conservation of elements, encoded by the linear constraints  $An = b$ , and non-negativity of molar quantities. In the ideal activity model,  $G$  is convex on  $\mathbb{R}_+^N$ , which ensures that (8) is a convex optimization problem. This classical formulation serves as the foundation for the existence and uniqueness analysis developed in Section 3.

## 3 On the minimizers of the Gibbs free energy

The current section gathers mathematical results on the optimization problem (8). We first establish in Section 3.1 a well-posedness result for the problem in close connection to the seminal work by Shapiro and Shapley [26]. Section 3.2 is then devoted to the thorough mathematical derivation of a unified formulation in the spirit of [19]. In such a formulation, a suitable notion of extended mole fractions is introduced for all the phases, including the absent ones. These persistent variables can then be used to solve the nonlinear system in an efficient way (see Section 4). We comment in Section 3.3 on the uniqueness of the solutions of the unified formulation. The loss of uniqueness due to azeotropy is described, as well as a possible lack of uniqueness for the extended mole fractions of the absent phases. The latter difficulty can however only happen for very specific right-hand sides, and uniqueness can be established as long as the two aforementioned identified difficulties do not occur.

### 3.1 Existence and uniqueness of the minimizers

The existence and uniqueness of a minimizer for an ideal multiphase chemical equilibrium have been rigorously examined by Shapiro and Shapley in [26]. In their fundamental contribution, they provide a comprehensive characterization of the minimizers set for the multiphase chemical equilibrium problem (8) in Theorem 9.2. Furthermore, they prove in Theorem 12.1 that if the configuration of the chemical system necessitates the presence of all phases at equilibrium, then the minimizer is unique. In cases where one or more phases are absent, uniqueness does not hold for the overall solution. However, Shapiro and Shapley demonstrate in Theorem 12.4 that the mole fractions of the present phases are uniquely defined. A last important result is Theorem 9.8, which states that if a species is present in a phase, then the entire phase must also be present.

In this section, we extend the work of Shapiro and Shapley by presenting a novel condition that ensures the uniqueness of a minimizer, even in cases where certain phases may be absent. This condition is derived from their characterization of the set of minimizers and provides a more comprehensive analysis framework.

Let  $\mathbb{R}_{>0}^N$  and  $\mathbb{R}_{\geq 0}^N$  be the set of positive and nonnegative vectors of  $\mathbb{R}^N$ , respectively. We define the set of vectors satisfying the constraints of conservation of elements by

$$\mathcal{M}_{\mathbf{A},\mathbf{b}} := \{\mathbf{n} \in \mathbb{R}^N \mid \mathbf{A}\mathbf{n} = \mathbf{b}\}.$$

We then make the following assumption:

$$(H1) \quad \mathcal{M}_{\mathbf{A},\mathbf{b}} \cap \mathbb{R}_{>0}^N \text{ is nonempty.}$$

This assumption corresponds to the existence of a strictly feasible point, which is known as the Slater condition in optimization. This assumption is physically reasonable, as it only requires that the element amounts  $\mathbf{b}$  do not a priori exclude any species from being present at equilibrium. Shapiro and Shapley have proved in [26, Theorem 9.2] that if the minimizers set for the multiphase chemical equilibrium problem (8) is nonempty, which is the case thanks to (H1) and since  $\mathcal{M}_{\mathbf{A},\mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$  is compact as shown in [18], then there exists a minimizer  $\tilde{\mathbf{n}}$  such that:

$$\arg \min \{G(\mathbf{m}) \mid \mathbf{m} \in \mathcal{M}_{\mathbf{A},\mathbf{b}} \cap \mathbb{R}_{\geq 0}^N\} = \mathcal{C}(\tilde{\mathbf{n}}) \cap \mathcal{M}_{\mathbf{A},\mathbf{b}} \cap \mathbb{R}_{\geq 0}^N, \quad (9)$$

where

$$\mathcal{C}(\tilde{\mathbf{n}}) = \{\mathbf{m} \in \mathbb{R}^N \mid \mathbf{m}^\alpha = \lambda_\alpha \tilde{\mathbf{n}}^\alpha, \lambda_\alpha \in \mathbb{R}, \forall \alpha = 1, \dots, N_{Ph}\}. \quad (10)$$

It follows from the definition of  $\mathcal{C}(\tilde{\mathbf{n}})$  that if  $\tilde{\mathbf{n}}^\alpha = \mathbf{0}$  then  $\mathbf{m}^\alpha = \mathbf{0}$ , for each  $\mathbf{m}$  belonging to  $\mathcal{C}(\tilde{\mathbf{n}}) \cap \mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$ . We introduce the set of present phases in  $\tilde{\mathbf{n}}$ , sometimes referred to as the *context* in the literature, by

$$\Gamma_{\tilde{\mathbf{n}}} := \{\alpha \in \{1, \dots, N_{Ph}\} \mid \tilde{\mathbf{n}}^\alpha \neq \mathbf{0}\}.$$

We can then establish the necessary and sufficient condition for the uniqueness of  $\tilde{\mathbf{n}}$  in the subsequent theorem.

**Theorem 3.1.** *Let  $\tilde{\mathbf{n}}$  the minimizer defined in (9) and (10) and let  $\mathbf{A}^\alpha$  be the  $M \times N^\alpha$  matrix corresponding to the part of  $\mathbf{A}$  related to the phase  $\alpha$ , then the minimizer  $\tilde{\mathbf{n}}$  is unique if and only if the set of vectors  $\{\mathbf{A}^\alpha \tilde{\mathbf{n}}^\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  is linearly independent.*

*Proof.* Let  $\mathbf{m} \in \mathcal{C}(\tilde{\mathbf{n}}) \cap \mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$  with  $\mathbf{m} \neq \tilde{\mathbf{n}}$ . Since both  $\tilde{\mathbf{n}}$  and  $\mathbf{m}$  belong to  $\mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$ , we have

$$\mathbf{A}(\tilde{\mathbf{n}} - \mathbf{m}) = \mathbf{0} \quad \Leftrightarrow \quad \sum_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}} \mathbf{A}^\alpha (\tilde{\mathbf{n}}^\alpha - \mathbf{m}^\alpha) = \mathbf{0}.$$

Since  $\mathbf{m} \in \mathcal{C}(\tilde{\mathbf{n}}) \cap \mathbb{R}_{\geq 0}^N$ , there exists  $\lambda_\alpha \geq 0$  such that  $\mathbf{m}^\alpha = \lambda_\alpha \tilde{\mathbf{n}}^\alpha$ ,  $\forall \alpha = 1, \dots, N_{Ph}$ . Therefore

$$\sum_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}} (1 - \lambda_\alpha) \mathbf{A}^\alpha \tilde{\mathbf{n}}^\alpha = \mathbf{0},$$

with  $\{1 - \lambda_\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  not all zero, meaning that  $\{\mathbf{A}^\alpha \tilde{\mathbf{n}}^\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  is not linearly independent.

Conversely, assume that  $\{\mathbf{A}^\alpha \tilde{\mathbf{n}}^\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  is not linearly independent, then there exists a set of non-zero scalars  $\{\gamma_\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  such that

$$\sum_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}} \gamma_\alpha \mathbf{A}^\alpha \tilde{\mathbf{n}}^\alpha = \mathbf{0}.$$

Since  $\tilde{\mathbf{n}} \geq \mathbf{0}$  and from (5), there exists  $\beta \in \Gamma_{\tilde{\mathbf{n}}}$  such that  $\max\{\gamma_\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}} =: \gamma_\beta > 0$ . We can define a set of non-negative scalars  $\{\lambda_\alpha\}_{\alpha \in \Gamma_{\tilde{\mathbf{n}}}}$  as  $\lambda_\alpha = 1 - \gamma_\alpha / \gamma_\beta$ . Therefore one defines the vector  $\mathbf{m}$  such that  $\mathbf{m}^\alpha = \lambda_\alpha \tilde{\mathbf{n}}^\alpha$  if  $\alpha \in \Gamma_{\tilde{\mathbf{n}}}$  and  $\mathbf{m}^\alpha = \mathbf{0}$  otherwise. By construction,  $\mathbf{m} \in \mathcal{C}(\tilde{\mathbf{n}}) \cap \mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$  and  $\mathbf{m} \neq \tilde{\mathbf{n}}$ .  $\square$

Theorem 3.1 extends the work of Shapiro and Shapley by guaranteeing the uniqueness of the total quantities in each phase, rather than merely the uniqueness of mole fractions. Specifically, it addresses a well-known issue of non-uniqueness related to azeotropic compositions. In such mixtures, the liquid and gaseous phases maintain the same composition at a fixed temperature corresponding to a constant boiling point, making it impossible to separate the components through simple distillation. This characteristic renders the two phases indistinguishable.

The notion of azeotropic mixture can be extended to our setting: we say that a minimizer  $\mathbf{n}$  of (8) has an azeotropic composition if  $\{\mathbf{A}^\alpha \mathbf{n}^\alpha\}_{\alpha \in \Gamma_{\mathbf{n}}}$  is linearly dependent, i.e. if

$$\text{there exists } (\lambda_\alpha)_{\alpha \in \Gamma_{\mathbf{n}}} \neq \mathbf{0} \text{ such that } \sum_{\alpha \in \Gamma_{\mathbf{n}}} \lambda_\alpha \mathbf{A}^\alpha \mathbf{n}^\alpha = \mathbf{0}. \quad (11)$$

We will further comment on the lack of uniqueness when (11) is satisfied in Section 3.3 later on.

Moreover, it serves as a clear example of the equality in mole fractions of two minimizers. The condition of Theorem 3.1 ensures that the total quantities of each phase in such cases are the same, thereby making the solution unique.

## 3.2 Unified formulation for multiphase chemical equilibria

In this section, we present our unified formulation for solving the minimization problem (8) using the concept of the subdifferential of a modified Gibbs energy that incorporates the non-negativity constraints. We establish the equivalence between this new set of equations and the original minimization problem, along with the uniqueness of the solution.

This approach also shows that the formulation of Coatléven and Michel [8] can be recovered directly from the KKT optimality conditions, thereby providing a rigorous theoretical basis for their method. Our formulation, however, adopts a different set of primary variables, offering an alternative perspective for solving the equilibrium equations.

This new formulation is presented in the following proposition.

**Proposition 3.1.** *Assume that hypothesis (H1) holds true. The vector  $\mathbf{n} \in \mathbb{R}^N$  is a minimizer for the problem (8) if and only if there exists a set of vectors  $(\boldsymbol{\xi}^\alpha, s_\alpha, r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  such that  $n_i = s_{\sigma(i)} \xi_i$ ,  $\xi_i > 0$ , for all  $i = 1, \dots, N$ , and satisfying:*

$$\sum_{\alpha=1}^{N_{Ph}} s_\alpha \mathbf{A}^\alpha \boldsymbol{\xi}^\alpha - \mathbf{b} = \mathbf{0}, \quad (12a)$$

$$\mathbf{S}^T [\boldsymbol{\mu}_i^\circ + \text{RT} \ln \xi_i]_{i=1, \dots, N} = \mathbf{0}, \quad (12b)$$

$$\langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle + r_\alpha - 1 = 0, \quad (\alpha = 1, \dots, N_{Ph}), \quad (12c)$$

$$s_\alpha r_\alpha = 0, \quad (\alpha = 1, \dots, N_{Ph}), \quad (12d)$$

$$s_\alpha \geq 0, r_\alpha \geq 0, \quad (\alpha = 1, \dots, N_{Ph}). \quad (12e)$$

In the formulation (12),  $\boldsymbol{\xi}^\alpha$  are the extended mole fractions, while  $s_\alpha$  represents the total quantities of phase  $\alpha$ . Consequently, if a phase is present at equilibrium, we have  $s_\alpha > 0$ , and  $\boldsymbol{\xi}^\alpha$  corresponds to the mole fractions:

$$\xi_i = \frac{n_i}{s_\alpha} = x_i(\mathbf{n}^\alpha).$$

Conversely, if a phase is absent,  $r_\alpha \geq 0$ , and the vector  $\boldsymbol{\xi}^\alpha$  does not necessarily sum to one:

$$\langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle \leq 1.$$

In this context,  $\boldsymbol{\xi}^\alpha$  serves as an extension of the mole fractions for vanishing phases, facilitating the computation of equilibria in a unified framework.

Note that the strict positivity of  $\boldsymbol{\xi}^\alpha$  implies that any solution to this system satisfies the result of [26, Theorem 9.8], which states that if a phase is present at equilibrium, then all species within that phase must also be present.

The remainder of this section will be on demonstrating Proposition 3.1.

*Proof of Proposition 3.1.* We introduce the sets of positive and non-negative vectors in  $\mathbb{R}^{N^\alpha}$  by  $\mathbb{R}_{>0}^{N^\alpha}$  and  $\mathbb{R}_{\geq 0}^{N^\alpha}$  respectively. Using these sets, we rewrite the minimization problem (8) as

$$\mathbf{n} \in \arg \min \left\{ \sum_{\alpha=1}^{N_{Ph}} \mathcal{G}_\alpha(\mathbf{n}^\alpha) \mid \mathbf{n} \in \mathcal{M}_{\mathbf{A}, \mathbf{b}} \right\}, \quad (13)$$

where

$$\mathcal{G}_\alpha(\mathbf{n}^\alpha) := \begin{cases} G_\alpha(\mathbf{n}^\alpha) & \text{if } \mathbf{n}^\alpha \in \mathbb{R}_{\geq 0}^{N^\alpha}, \\ +\infty & \text{otherwise.} \end{cases} \quad (14)$$

In this way, non-negativity constraints are incorporated into the objective function.

The Gibbs function  $\mathcal{G}_\alpha$  in (14) is only differentiable in  $\mathbb{R}_{>0}^{N^\alpha}$ , but it is possible to define its subdifferential for any  $\mathbf{n}^\alpha \in \mathbb{R}^{N^\alpha}$ . The subdifferential of the convex function  $\mathcal{G}_\alpha : \mathbb{R}^{N^\alpha} \rightarrow \mathbb{R}$  at a point  $\mathbf{n}^\alpha \in \mathbb{R}^{N^\alpha}$  is denoted  $\partial \mathcal{G}_\alpha(\mathbf{n}^\alpha)$  and is defined by [25, Section 23]:

$$\boldsymbol{\mu}^\alpha \in \partial \mathcal{G}_\alpha(\mathbf{n}^\alpha) \iff \mathcal{G}_\alpha(\mathbf{m}^\alpha) \geq \mathcal{G}_\alpha(\mathbf{n}^\alpha) + \langle \boldsymbol{\mu}^\alpha, \mathbf{m}^\alpha - \mathbf{n}^\alpha \rangle, \quad \forall \mathbf{m}^\alpha \in \mathbb{R}^{N^\alpha}.$$

The function  $\boldsymbol{\mu}^\alpha$  is called a subgradient of  $\mathcal{G}_\alpha$ . In our particular case, one has

$$\partial \mathcal{G}_\alpha(\mathbf{n}^\alpha) = \begin{cases} \{\nabla \mathcal{G}_\alpha(\mathbf{n}^\alpha)\} & \text{if } \mathbf{n}^\alpha \in \mathbb{R}_{>0}^{N^\alpha}, \\ \emptyset & \text{otherwise.} \end{cases}$$

For  $\mathbf{n}^\alpha = \mathbf{0}$ , the subdifferential  $\partial \mathcal{G}_\alpha(\mathbf{0})$  contains infinitely many elements, as shows the analysis carried out in what follows, cf. Lemma 3.1.

Then, if  $\mathbf{n} = (\mathbf{n}^\alpha)_{\alpha=1, \dots, N_{Ph}}$  is a solution to the problem defined in (13), it must satisfy the first-order KKT optimality conditions:

$$\mathbf{A} \mathbf{n} - \mathbf{b} = \mathbf{0}, \quad (15a)$$

$$\boldsymbol{\mu} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0}, \quad (15b)$$

$$\boldsymbol{\mu} = (\boldsymbol{\mu}^\alpha)_{\alpha=1, \dots, N_{Ph}} \quad (15c)$$

$$\boldsymbol{\mu}^\alpha \in \partial\mathcal{G}_\alpha(\mathbf{n}^\alpha). \quad (15d)$$

In (15), the subdifferential of  $\mathcal{G}_\alpha$  includes the non-negativity constraint, while  $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_M)^T$  represents the vector of Lagrange multipliers associated with the linear constraint (15a).

To establish our formulation, we need to characterize  $\partial\mathcal{G}_\alpha(\mathbf{n}^\alpha)$ . The case where  $\mathbf{n}^\alpha < \mathbf{0}$  is excluded from the minimization problem (13) by the definition of  $\mathcal{G}_\alpha$ . The following lemma provides a characterization of  $\partial\mathcal{G}_\alpha(\mathbf{n}^\alpha)$  for  $\mathbf{n}^\alpha \geq \mathbf{0}$ .

**Lemma 3.1.** *Let  $\mathbf{n}^\alpha \geq \mathbf{0}$ , then  $\boldsymbol{\mu}^\alpha \in \partial\mathcal{G}_\alpha(\mathbf{n}^\alpha)$  if and only if there exists  $s_\alpha, \boldsymbol{\xi}^\alpha = (\xi_i)_{i \in \sigma^{-1}(\alpha)}$  such that:*

$$s_\alpha \geq 0, \boldsymbol{\xi}^\alpha > \mathbf{0}, \quad \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle \leq 1, \quad s_\alpha (1 - \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle) = 0,$$

and for all  $i \in \sigma^{-1}(\alpha)$ ,

$$\mu_i = \mu_i^0 + RT \ln \xi_i, \quad n_i = s_\alpha \xi_i.$$

*Proof of Lemma 3.1.* In order to prove Lemma 3.1, we will describe the subdifferential of the Legendre transform  $\mathcal{G}_\alpha^*$  of  $\mathcal{G}_\alpha$ . Indeed, since  $\mathcal{G}_\alpha$  is a lower semi-continuous proper convex function, the following equivalence holds:

$$\boldsymbol{\mu}^\alpha \in \partial\mathcal{G}_\alpha(\mathbf{n}^\alpha) \Leftrightarrow \mathbf{n}^\alpha \in \partial\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha), \quad (16)$$

where  $\mathcal{G}_\alpha^*$  is the Legendre transform of  $\mathcal{G}_\alpha$ . This transformation is widely used in thermodynamics to convert a function defined in one set of variables into another, effectively representing the function in terms of its tangents. It is defined as

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = \sup_{\mathbf{n}^\alpha \in \mathbb{R}^{N_\alpha}} \sum_{i \in \sigma^{-1}(\alpha)} n_i \mu_i - \mathcal{G}_\alpha(\mathbf{n}^\alpha), \quad \boldsymbol{\mu}^\alpha \in \mathbb{R}^{N_\alpha}. \quad (17)$$

For  $\mathbf{n}^\alpha \notin \mathbb{R}_{\geq 0}^{N_\alpha}$ , we have  $\mathcal{G}_\alpha(\mathbf{n}^\alpha) = +\infty$ , we can then restrict the supremum in (17) to  $\mathbf{n}^\alpha \geq \mathbf{0}$  and consider the function  $G_\alpha$  instead of  $\mathcal{G}_\alpha$ . We thus obtain

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = \sup_{\mathbf{n}^\alpha \geq \mathbf{0}} \sum_{i \in \sigma^{-1}(\alpha)} n_i \mu_i - G_\alpha(\mathbf{n}^\alpha) = \sup_{\mathbf{n}^\alpha \geq \mathbf{0}} \sum_{i \in \sigma^{-1}(\alpha)} n_i [\mu_i - \mu_i^\circ - RT \ln x_i(\mathbf{n}^\alpha)].$$

Shifting from the supremum to the infimum, we get

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = - \inf_{\mathbf{n}^\alpha \geq \mathbf{0}} \sum_{i \in \sigma^{-1}(\alpha)} n_i [\mu_i^\circ + RT \ln x_i(\mathbf{n}^\alpha) - \mu_i]. \quad (18)$$

We then introduce mole fractions by means of the following change of variables:

$$x_i s_\alpha = n_i \quad \text{with} \quad s_\alpha = \langle \mathbf{n}^\alpha, \mathbf{1} \rangle.$$

The Legendre transform (18) is then written as

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = - \inf_{\substack{s_\alpha, \mathbf{x} \geq \mathbf{0} \\ \langle \mathbf{x}^\alpha, \mathbf{1} \rangle = 1}} s_\alpha \sum_{i \in \sigma^{-1}(\alpha)} x_i [\mu_i^\circ + RT \ln x_i - \mu_i]. \quad (19)$$

From this relation, we can extract a minimization problem by defining the convex function  $g_\alpha(\cdot; \boldsymbol{\mu}^\alpha) : \mathbb{R}_{>0}^{N_\alpha} \rightarrow \mathbb{R}$  defined as

$$g_\alpha(\mathbf{x}^\alpha; \boldsymbol{\mu}^\alpha) = \sum_{i \in \sigma^{-1}(\alpha)} x_i [\mu_i^\circ + RT \ln x_i - \mu_i], \quad (20)$$

where  $\boldsymbol{\mu}^\alpha$  is considered as a parameter. Then, introducing the convex sets

$$X_\alpha = \{\mathbf{x}^\alpha > \mathbf{0} \mid \langle \mathbf{x}^\alpha, \mathbf{1} \rangle = 1\}, \quad \text{and} \quad \overline{X}_\alpha = \{\mathbf{x}^\alpha \geq \mathbf{0} \mid \langle \mathbf{x}^\alpha, \mathbf{1} \rangle = 1\},$$

the relation (19) becomes

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = - \inf_{s_\alpha \geq 0} s_\alpha \min_{\mathbf{x}^\alpha \in \overline{X}_\alpha} g_\alpha(\mathbf{x}^\alpha; \boldsymbol{\mu}^\alpha). \quad (21)$$

We can then study the minimization problem:

$$\min_{\mathbf{x}^\alpha \in \overline{X}_\alpha} g_\alpha(\mathbf{x}^\alpha; \boldsymbol{\mu}^\alpha). \quad (22)$$

The following lemma, to be proven at the end of the subsection, states that the inequality constraint on  $\mathbf{x}^\alpha$  is never saturated at the solution of (22).

**Lemma 3.2.** *If  $\mathbf{x}^* \in \overline{X_\alpha}$  minimizes  $g_\alpha(\cdot; \boldsymbol{\mu}^\alpha)$  on  $\overline{X_\alpha}$ , then  $\mathbf{x}^* \in X_\alpha$ .*

Based on Lemma 3.2, the problem (22) can be simplified to

$$\min_{\langle \mathbf{x}^\alpha, \mathbf{1} \rangle - 1 = 0} \sum_{i \in \sigma^{-1}(\alpha)} x_i [\mu_i^\circ + \text{RT} \ln x_i - \mu_i]. \quad (23)$$

The Euler-Lagrange equations associated to the problem (23) are

$$\mu_i^\circ + \text{RT} \ln x_i - \mu_i + \text{RT} - \lambda = 0, \quad i \in \sigma^{-1}(\alpha), \quad (24)$$

$$\langle \mathbf{x}^\alpha, \mathbf{1} \rangle - 1 = 0, \quad (25)$$

where  $\lambda$  is the Lagrange multiplier for the constraint (25). Applying equation (24) and then equation (25) to the problem (23) yields the following minimum:

$$\min_{\langle \mathbf{x}^\alpha, \mathbf{1} \rangle - 1 = 0} \sum_{i \in \sigma^{-1}(\alpha)} x_i [\mu_i^\circ + \text{RT} \ln x_i - \mu_i] = \min_{\langle \mathbf{x}^\alpha, \mathbf{1} \rangle - 1 = 0} \sum_{i \in \sigma^{-1}(\alpha)} x_i [\lambda - \text{RT}] = \lambda - \text{RT}.$$

In order to introduce this minimum into the Legendre transform (21), we will rewrite it in terms of  $\boldsymbol{\mu}^\alpha$ . To do so, we define the functions  $\xi_i : \mathbb{R} \rightarrow \mathbb{R}_+^*$ ,  $i \in \sigma^{-1}(\alpha)$ , by

$$\xi_i(\mu_i) := \exp\left(\frac{\mu_i - \mu_i^\circ}{\text{RT}}\right). \quad (26)$$

Then using (24) and (25), we get that:

$$\xi_i(\mu_i) = x_i \exp\left(\frac{\text{RT} - \lambda}{\text{RT}}\right) \quad \text{and} \quad \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle = \sum_{i \in \sigma^{-1}(\alpha)} x_i \exp\left(\frac{\text{RT} - \lambda}{\text{RT}}\right) = \exp\left(\frac{\text{RT} - \lambda}{\text{RT}}\right),$$

where  $\boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha) = (\xi_i(\mu_i))_{i \in \sigma^{-1}(\alpha)}$ . The solution of problem (22) is then given by:

$$\min_{\mathbf{x}^\alpha \in \overline{X_\alpha}} g_\alpha(\mathbf{x}^\alpha) = \lambda - \text{RT} = -\text{RT} \ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle. \quad (27)$$

Using the minimum (27) into (21), we obtain:

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = - \inf_{s_\alpha \geq 0} -s_\alpha \text{RT} \ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle = \sup_{s_\alpha \geq 0} s_\alpha \text{RT} \ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle.$$

Therefore, since  $s_\alpha \text{RT} \geq 0$ , we can consider two cases based on the sign of  $\ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle$ :

- $\langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle \leq 1 \quad \Rightarrow \quad \ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle \leq 0 \quad \Rightarrow \quad \mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = 0;$
- $\langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle > 1 \quad \Rightarrow \quad \ln \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle > 0 \quad \Rightarrow \quad \mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = +\infty.$

We define the set  $K_\alpha$  as follows

$$K_\alpha := \{\boldsymbol{\mu}^\alpha \mid \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle \leq 1\}.$$

This set is a convex set as it is the sublevel set of the convex function  $\Psi : \boldsymbol{\mu}^\alpha \mapsto \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle$ . Alongside  $K_\alpha$ , we introduce the associated characteristic function:

$$\chi_{K_\alpha}(\boldsymbol{\mu}^\alpha) = \begin{cases} 0 & \text{if } \boldsymbol{\mu}^\alpha \in K_\alpha, \\ +\infty & \text{if } \boldsymbol{\mu}^\alpha \notin K_\alpha. \end{cases}$$

With these definitions, we can express the Legendre transform as follows:

$$\mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) = \chi_{K_\alpha}(\boldsymbol{\mu}^\alpha).$$

We can now easily calculate the subdifferential of  $\mathcal{G}_\alpha^*$  at a point  $\boldsymbol{\mu}^\alpha \in K_\alpha$ . If  $\boldsymbol{\mu}^\alpha$  lies in the interior of  $K_\alpha$ , the subdifferential is reduced to the gradient of  $\mathcal{G}_\alpha^*$ , which is equal to zero. Otherwise, if  $\boldsymbol{\mu}^\alpha$  is on the boundary of  $K_\alpha$ , the definition of the subdifferential leads to the condition:

$$\mathbf{n}^\alpha \in \partial \mathcal{G}_\alpha^*(\boldsymbol{\mu}^\alpha) \quad \Leftrightarrow \quad 0 \geq \langle \mathbf{n}^\alpha, \boldsymbol{\eta}^\alpha - \boldsymbol{\mu}^\alpha \rangle, \quad \forall \boldsymbol{\eta}^\alpha \in \mathbb{R}^{N^\alpha},$$

which is known as the normal cone of the convex set  $K_\alpha$ . The boundary of  $K_\alpha$  is smooth because it is defined by  $\Psi(\boldsymbol{\mu}^\alpha) = 1$ . Consequently, this cone can be simplified to the normal direction, which is determined by the gradient of  $\Psi$ . Therefore, the subdifferential of  $\mathcal{G}_\alpha^\star$  is:

$$\partial\mathcal{G}_\alpha^\star(\boldsymbol{\mu}^\alpha) = \begin{cases} \mathbf{0} & \text{if } \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle < 1, \\ \gamma_\alpha \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \gamma_\alpha \geq 0 & \text{if } \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle = 1. \end{cases}$$

We can now conclude the proof of Lemma 3.1.

Let  $\mathbf{n}^\alpha \geq \mathbf{0}$ , then

$$\begin{aligned} \boldsymbol{\mu}^\alpha \in \partial\mathcal{G}_\alpha(\mathbf{n}^\alpha) &\Leftrightarrow \mathbf{n}^\alpha \in \partial\mathcal{G}_\alpha^\star(\boldsymbol{\mu}^\alpha) \Leftrightarrow \begin{cases} n_i = \gamma_\alpha \xi_i(\mu_i), \gamma_\alpha \geq 0, i \in \sigma^{-1}(\alpha), \\ \langle \boldsymbol{\xi}^\alpha(\boldsymbol{\mu}^\alpha), \mathbf{1} \rangle \leq 1, \end{cases} \\ &\Leftrightarrow \begin{cases} n_i = \gamma_\alpha \xi_i, \gamma_\alpha \geq 0, i \in \sigma^{-1}(\alpha), \\ \mu_i = \mu_i^\circ + \text{RT} \ln \xi_i, \xi_i > 0, \\ \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle \leq 1. \end{cases} \end{aligned}$$

We already know that if  $\mathbf{n}^\alpha > \mathbf{0}$ , the derivative of  $g_\alpha$  is given by  $\mu_i(\mathbf{n}^\alpha) = \mu_i^\circ + \text{RT} \ln n_i/s_\alpha$ , thus  $\xi_i = n_i/s_\alpha$ ,  $\gamma_\alpha = s_\alpha > 0$  and  $\langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle = 1$ . In the case where  $\mathbf{n}^\alpha = \mathbf{0}$ , it is clear that  $\gamma_\alpha = 0 = s_\alpha$  and that  $\langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle \leq 1$ . Therefore, in each of these cases, one has

$$s_\alpha \geq 0, \boldsymbol{\xi}^\alpha > \mathbf{0}, \quad \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle \leq 1, \quad s_\alpha (1 - \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle) = 0, \quad n_i = s_\alpha \xi_i, \quad i \in \sigma^{-1}(\alpha),$$

concluding the proof of Lemma 3.1.  $\square$

Using this result, we can rewrite the KKT optimality conditions (15) as

$$\begin{aligned} \sum_{\alpha=1}^{N_{Ph}} s_\alpha \mathbf{A}^\alpha \boldsymbol{\xi}^\alpha - \mathbf{b} &= \mathbf{0}, \\ \mathbf{S}^T [\mu_i^\circ + \text{RT} \ln \xi_i]_{i=1, \dots, N} &= \mathbf{0}, \\ s_\alpha (1 - \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle) &= 0, & (\alpha = 1, \dots, N_{Ph}), \\ s_\alpha \geq 0, 1 - \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle &\geq 0, & (\alpha = 1, \dots, N_{Ph}). \end{aligned}$$

where  $\boldsymbol{\mu}(\boldsymbol{\xi}) := (\mu_i^\circ + \text{RT} \ln \xi_i)_{i=1, \dots, N}$ . Therefore, based on the relation  $\mathbf{S}^T \mathbf{A}^T = \mathbf{0}$  from (4), and defining

$$r_\alpha := 1 - \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle,$$

we finally obtain:

$$\begin{aligned} \sum_{\alpha=1}^{N_{Ph}} s_\alpha \mathbf{A}^\alpha \boldsymbol{\xi}^\alpha - \mathbf{b} &= \mathbf{0}, & (12a) \\ \mathbf{S}^T [\mu_i^\circ + \text{RT} \ln \xi_i]_{i=1, \dots, N} &= \mathbf{0}, & (12b) \\ \langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle + r_\alpha - 1 &= 0, & (\alpha = 1, \dots, N_{Ph}), & (12c) \\ s_\alpha r_\alpha &= 0, & (\alpha = 1, \dots, N_{Ph}), & (12d) \\ s_\alpha \geq 0, r_\alpha &\geq 0, & (\alpha = 1, \dots, N_{Ph}). & (12e) \end{aligned}$$

Given the convexity of the Gibbs function and constraints, the equivalence between the minimization problem and the solution of system (12) follows from assumption (H1), which guarantees that the Slater condition holds – namely, the existence of a strictly feasible point. This classical constraint qualification ensures strong duality and the equivalence between the constrained minimization problem and its KKT system. This concludes the proof of Proposition 3.1.  $\square$

It remains to proof Lemma 3.2 to complete the proof of the main result of the section, namely Proposition 3.1.

*Proof of Lemma 3.2.* The function  $g_\alpha$  is continuous on the compact set  $\overline{X_\alpha}$ , then from the Weierstrass theorem, there exists  $\mathbf{x}^\star \in \{\arg \min g_\alpha(\mathbf{x}^\alpha; \boldsymbol{\mu}^\alpha) \mid \mathbf{x}^\alpha \in \overline{X_\alpha}\}$ . Let us assume that  $\mathbf{x}^\star \in \overline{X_\alpha} \setminus X_\alpha$ , meaning that there exists a subset  $\mathcal{J} \subsetneq \sigma^{-1}(\alpha)$  such that  $x_j^\star = 0$ , for all  $j \in \mathcal{J}$ . Note that the case  $\mathcal{J} = \sigma^{-1}(\alpha)$  is excluded since  $\langle \boldsymbol{\xi}^\alpha, \mathbf{1} \rangle = 1$ . Let  $\mathbf{x} \in X_\alpha$  and  $\varepsilon \in (0, 1)$ , then one defines  $\mathbf{x}^0 := \mathbf{x} - \mathbf{x}^\star$  and  $\mathbf{x}^\varepsilon := \mathbf{x}^\star + \varepsilon \mathbf{x}^0 = \mathbf{x}^\varepsilon = \varepsilon \mathbf{x} + (1 - \varepsilon) \mathbf{x}^\star$ . The vector  $\mathbf{x}^\varepsilon$  is a convex linear combination of vectors of  $\overline{X_\alpha}$

which is a convex set, hence  $\mathbf{x}^\varepsilon \in \overline{X_\alpha}$ . Furthermore,  $\mathbf{x}^\varepsilon \in X_\alpha$  since  $\mathbf{x}^\varepsilon = \varepsilon \mathbf{x} + (1 - \varepsilon) \mathbf{x}^* \geq \varepsilon \mathbf{x} > 0$ . By convexity of  $g_\alpha$  on  $\mathbb{R}_{\geq 0}^{N_\alpha}$ ,

$$g_\alpha(\mathbf{x}^*; \boldsymbol{\mu}^\alpha) \geq g_\alpha(\mathbf{x}^\varepsilon; \boldsymbol{\mu}^\alpha) + \langle \nabla g_\alpha(\mathbf{x}^\varepsilon), \mathbf{x}^* - \mathbf{x}^\varepsilon \rangle \quad (28)$$

$$\Leftrightarrow \frac{g_\alpha(\mathbf{x}^*; \boldsymbol{\mu}^\alpha) - g_\alpha(\mathbf{x}^\varepsilon; \boldsymbol{\mu}^\alpha)}{\varepsilon} \geq -\langle \nabla g_\alpha(\mathbf{x}^\varepsilon), \mathbf{x}^0 \rangle, \quad (29)$$

where  $\nabla g_\alpha(\mathbf{x}^\varepsilon) := (\mu_i^\circ + RT \ln x_i^\varepsilon + RT)_{i \in \sigma^{-1}(\alpha)}$ .

We will now take the limit when  $\varepsilon$  tends to 0 in inequality (29). In the right-hand side one has

$$\lim_{\varepsilon \rightarrow 0} -\langle \nabla g_\alpha(\mathbf{x}^\varepsilon), \mathbf{x}^0 \rangle = - \sum_{j \in \mathcal{J}} x_j^0 \lim_{\varepsilon \rightarrow 0} \partial_{x_j} g_\alpha(\mathbf{x}^\varepsilon) - \sum_{i \in \sigma^{-1}(\alpha) \setminus \mathcal{J}} x_i^0 \lim_{\varepsilon \rightarrow 0} \partial_{x_i} g_\alpha(\mathbf{x}^\varepsilon). \quad (30)$$

Noting that  $\lim_{\varepsilon \rightarrow 0} \mathbf{x}^\varepsilon = \mathbf{x}^*$  and in particular that  $\lim_{\varepsilon \rightarrow 0} x_j^\varepsilon = x_j^* = 0$ , for all  $j \in \mathcal{J}$ , it follows from the continuity of  $\nabla g_\alpha$  that

$$\lim_{\varepsilon \rightarrow 0} \partial_{x_j} g_\alpha(\mathbf{x}^\varepsilon) = -\infty, \forall j \in \mathcal{J} \quad \text{and} \quad \lim_{\varepsilon \rightarrow 0} \partial_{x_i} g_\alpha(\mathbf{x}^\varepsilon) = \partial_{x_i} g_\alpha(\mathbf{x}^*) \in \mathbb{R}, \forall i \in \sigma^{-1}(\alpha) \setminus \mathcal{J}. \quad (31)$$

By combining (30) and (31) with  $x_j^0 = x_j > 0$ , for all  $j \in \mathcal{J}$ , one finds that the right-hand side of (29) tends to  $+\infty$ . However if  $\mathbf{x}^*$  minimizes  $g_\alpha$  on  $\overline{X_\alpha}$ , then the left-hand side of (29) is non-positive which is a contradiction. Therefore  $\mathcal{J} = \emptyset$  and  $\mathbf{x}^* \in X_\alpha$ .  $\square$

### 3.3 On the uniqueness for the unified formulation

The existence of a solution is guaranteed by the assumption that  $\mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{> 0}^N$  is nonempty (H1) and by the compactness of  $\mathcal{M}_{\mathbf{A}, \mathbf{b}} \cap \mathbb{R}_{\geq 0}^N$ . We have seen in Theorem 3.1 that the minimizer  $\mathbf{n}$  for (8), defined in (9) and (10), is unique if and only if the set of vectors  $\{\mathbf{A}^\alpha \mathbf{n}^\alpha\}_{\alpha \in \Gamma_{\mathbf{n}}}$ , where  $\Gamma_{\mathbf{n}} = \{\alpha \in \{1, \dots, N_{Ph}\} \mid \mathbf{n}^\alpha \neq \mathbf{0}\}$ , is linearly independent.

To better understand this condition, let us assume that condition (11) holds true for some minimizer  $\mathbf{n} = (s_{\sigma(i)} \xi_i)_{i \in \mathcal{C}}$  of the Gibbs energy, with  $s_\alpha$  and  $\boldsymbol{\xi}^\alpha$ ,  $\alpha \in \mathcal{P}$ , as in Proposition 3.1. Then setting  $\lambda_\alpha = 0$  for all  $\alpha \in \mathcal{P} \setminus \Gamma_{\mathbf{n}}$ , then  $(\lambda_\alpha)_{\alpha \in \mathcal{P}} \neq \mathbf{0}$  is such that

$$\sum_{\alpha \in \mathcal{P}} \lambda_\alpha \mathbf{A}^\alpha \boldsymbol{\xi}^\alpha = \mathbf{0},$$

or, equivalently, such that

$$\mathbf{w} = (\lambda_{\sigma(i)} \xi_i)_{i \in \mathcal{C}} \in \ker \mathbf{A} \setminus \mathbf{0}. \quad (32)$$

Since  $\lambda_\alpha = 0$  if  $s_\alpha = 0$ , there exists  $\bar{\varepsilon} > 0$  such that  $s_\alpha + \varepsilon \lambda_\alpha \geq 0$  for all  $\varepsilon \in (0, \bar{\varepsilon})$ . Denote by

$$\mathbf{n}_\varepsilon = ((s_{\sigma(i)} + \varepsilon \lambda_{\sigma(i)}) \xi_i)_{i \in \mathcal{C}} \geq \mathbf{0},$$

then  $\mathbf{A} \mathbf{n}_\varepsilon = \mathbf{b}$ , and, setting  $\boldsymbol{\mu} = (\mu_i^\circ + RT \log(\xi_i))_{i \in \mathcal{C}} = (\mu_i)_{i \in \mathcal{C}}$  as in the proof of Proposition 3.1, we check that

$$G(\mathbf{n}_\varepsilon) = \sum_{\alpha \in \mathcal{P}} (s_\alpha + \varepsilon \lambda_\alpha) \sum_{i \in \sigma^{-1}(\alpha)} \xi_i \mu_i = G(\mathbf{n}) + \varepsilon \mathbf{w}^T \boldsymbol{\mu} = G(\mathbf{n})$$

since  $\boldsymbol{\mu} \in \ker \mathbf{S}^T$  owing to (12b), since  $\mathbf{w} \in \ker \mathbf{A}$  owing to (32), and since  $\ker \mathbf{S}^T = (\ker \mathbf{A})^\perp$  thanks to (4). Therefore,  $\mathbf{n}_\varepsilon$  is also minimizer of problem (8), thus a solution to the unified formulation of Proposition 3.1, for all  $\varepsilon \in (0, \bar{\varepsilon})$ . In particular, condition (11) yields an infinite number of solutions to the minimization problem, and thus to the unified formulation.

While condition

$$(\mathbf{A}^\alpha \mathbf{n}^\alpha)_{\alpha \in \Gamma_{\mathbf{n}}} \text{ is linearly independent}$$

is intrinsic to the physical problem to avoid azeotropy, and does not stem from our unified formulation, uniqueness can also fail in some very specific configurations where the extended mole fraction may become ambiguous for some absent species. To illustrate such configurations, let us analyze the following example system:

$$\begin{aligned} \mathcal{C} &= \{\text{H}_2\text{O}(\text{aq}), \text{H}^+, \text{OH}^-, \text{H}_2\text{O}(\text{g})\}, \\ \mathcal{E} &= \{\text{H}, \text{O}\}, \\ \mathcal{R} &= \{ \text{OH}^- = \text{H}_2\text{O}(\text{aq}) - \text{H}^+, \\ &\quad \text{H}_2\text{O}(\text{g}) = \text{H}_2\text{O}(\text{aq}) \}, \\ \mathcal{P} &= \{\text{aqueous}, \text{gaseous}\}, \end{aligned} \quad (33)$$

with the right-hand side

$$\mathbf{b} = \lambda(2, 1)^T, \lambda > 0,$$

and with given pressure and temperature such that  $\mu_{\text{H}_2\text{O}(\text{g})}^\circ < \mu_{\text{H}_2\text{O}(\text{aq})}^\circ$ . The associated formula and stoichiometry matrices are

$$\mathbf{A} = \begin{pmatrix} 2 & 1 & 1 & 2 \\ 1 & 0 & 1 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{S}^T = \begin{pmatrix} 1 & -1 & -1 & 0 \\ 1 & 0 & 0 & -1 \end{pmatrix}.$$

Let

$$\boldsymbol{\xi} = (\xi_{\text{H}_2\text{O}(\text{aq})}, \xi_{\text{H}^+}, \xi_{\text{OH}^-}, \xi_{\text{H}_2\text{O}(\text{g})}), \quad \mathbf{s} = (s_{\text{aq}}, s_{\text{g}}) \quad \text{and} \quad \mathbf{r} = (r_{\text{aq}}, r_{\text{g}})$$

be such that

$$\xi_{\text{H}_2\text{O}(\text{g})} = 1 \quad \text{and} \quad s_{\text{aq}} = 0.$$

Let us establish under which conditions the system (12) is fulfilled. Equation (12a) yields

$$s_{\text{aq}} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \xi_{\text{H}_2\text{O}(\text{aq})} \\ \xi_{\text{H}^+} \\ \xi_{\text{OH}^-} \end{pmatrix} + s_{\text{g}} \begin{pmatrix} 2 \\ 1 \end{pmatrix} \xi_{\text{H}_2\text{O}(\text{g})} = \lambda \begin{pmatrix} 2 \\ 1 \end{pmatrix},$$

and is satisfied if  $s_{\text{g}} = \lambda$ .

It follows from equations (12c) and (12d) that  $(\xi_{\text{H}_2\text{O}(\text{aq})}, \xi_{\text{H}^+}, \xi_{\text{OH}^-})$  must satisfy

$$\xi_{\text{H}_2\text{O}(\text{aq})} + \xi_{\text{H}^+} + \xi_{\text{OH}^-} \leq 1 \tag{34}$$

and that  $r_{\text{g}} = 0, r_{\text{aq}} \geq 0$ . Finally, equations (12b) yields

$$\frac{\mu_{\text{H}_2\text{O}(\text{aq})}^\circ - \mu_{\text{H}^+}^\circ - \mu_{\text{OH}^-}^\circ}{RT} + \ln \xi_{\text{H}_2\text{O}(\text{aq})} - \ln \xi_{\text{H}^+} - \ln \xi_{\text{OH}^-} = 0, \tag{35a}$$

$$\frac{\mu_{\text{H}_2\text{O}(\text{aq})}^\circ - \mu_{\text{H}_2\text{O}(\text{g})}^\circ}{RT} + \ln \xi_{\text{H}_2\text{O}(\text{aq})} - \underbrace{\ln \xi_{\text{H}_2\text{O}(\text{g})}}_{=0} = 0. \tag{35b}$$

We can rewrite (35b) as

$$\xi_{\text{H}_2\text{O}(\text{aq})} = \underbrace{e^{\frac{\mu_{\text{H}_2\text{O}(\text{g})}^\circ - \mu_{\text{H}_2\text{O}(\text{aq})}^\circ}{RT}}}_{:=\kappa_1}. \tag{36}$$

Therefore, using (36) into (35a), we get that

$$\xi_{\text{OH}^-} = \frac{1}{\xi_{\text{H}^+}} \underbrace{e^{\frac{\mu_{\text{H}_2\text{O}(\text{g})}^\circ - \mu_{\text{H}^+}^\circ - \mu_{\text{OH}^-}^\circ}{RT}}}_{:=\kappa_2}. \tag{37}$$

It follows from (34), (36) and (37) that:

$$\xi_{\text{H}_2\text{O}(\text{aq})} + \xi_{\text{H}^+} + \xi_{\text{OH}^-} \leq 1 \Leftrightarrow \xi_{\text{H}^+}^2 + (\kappa_1 - 1)\xi_{\text{H}^+} + \kappa_2 \leq 0.$$

The discriminant of this quadratic inequality is given by:

$$\Delta = (\kappa_1 - 1)^2 - 4\kappa_2.$$

Using (34), (36) and (37) we get that:

$$1 - \kappa_1 \geq \xi_{\text{H}^+} + \xi_{\text{OH}^-} \Leftrightarrow \Delta \geq (\xi_{\text{H}^+} + \xi_{\text{OH}^-})^2 - 4\xi_{\text{H}^+}\xi_{\text{OH}^-} \Leftrightarrow \Delta \geq (\xi_{\text{H}^+} - \xi_{\text{OH}^-})^2 \geq 0.$$

It follows that:

$$\xi_{\text{H}^+} \in \left[ \frac{1 - \kappa_1 - \sqrt{(1 - \kappa_1)^2 - 4\kappa_2}}{2}, \frac{1 - \kappa_1 + \sqrt{(1 - \kappa_1)^2 - 4\kappa_2}}{2} \right].$$

Therefore, the chemical system (33) leads to a nonunique solution of equations (12). Roughly speaking, in a configuration with only pure water vapor, our unified formulation cannot give a precise value for the pH of the absent liquid phase.

The previous example indicates that a condition on the vector  $\mathbf{b}$  is necessary to ensure the uniqueness of the vector  $\boldsymbol{\xi}$ . This condition basically requires that the vector  $\mathbf{b}$  cannot be recovered with less than  $M$  species. Recalling that  $\mathbf{a}_j$  is the formula vector for the species  $C_j$  as defined in Section 2, let us make the following assumptions about the vector  $\mathbf{b}$ :

(H2) for all subset  $\Gamma \subset \{1, \dots, N_{Ph}\}$  such that  $\{\mathbf{a}_j\}_{j \in \sigma^{-1}(\Gamma)}$  is a non-spanning set of vectors in  $\mathbb{R}^M$ , one has  $\mathbf{b} \notin \text{span}(\{\mathbf{a}_j\}_{j \in \sigma^{-1}(\Gamma)})$ .

We can then write the following uniqueness proposition.

**Proposition 3.2.** *Let  $\mathcal{X} = (\xi^\alpha, s_\alpha, r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  be a solution of system (12) where  $\mathbf{b}$  satisfies assumption (H2), and let  $\Gamma_{\mathcal{X}}$  be the corresponding set of present phases. Then  $\mathcal{X}$  is unique if and only if the set of vectors  $\{\mathbf{A}^\alpha \xi^\alpha\}_{\alpha \in \Gamma_{\mathcal{X}}}$  is linearly independent.*

To prove this proposition, we first need to establish the following two intermediate results.

**Lemma 3.3.** *Let  $\mathcal{X} = (\xi^\alpha, s_\alpha, r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  be a solution of (12). If  $\mathbf{b}$  satisfies assumption (H2), then there are at least  $M$  species present in  $\mathcal{X}$ . Furthermore, system (12) can be rewritten with these  $M$  species as the primary species.*

*Proof.* Let  $\Gamma_{\mathcal{X}}$  be the set of present phases of  $\mathcal{X}$ , in this way  $\sigma^{-1}(\Gamma_{\mathcal{X}})$  is the set of all present species. From the assumption (H2),  $\{\mathbf{a}_j\}_{j \in \sigma^{-1}(\Gamma_{\mathcal{X}})}$  is a spanning set in  $\mathbb{R}^M$ . It is then possible to extract a basis of  $\mathbb{R}^M$  from the set  $\{\mathbf{a}_j\}_{j \in \sigma^{-1}(\Gamma_{\mathcal{X}})}$  and to construct a permutation matrix  $\mathbf{P} \in \mathbb{R}^{N \times N}$  such that this basis forms a primary block in  $\tilde{\mathbf{A}} := \mathbf{A}\mathbf{P}$ . In this way,  $\tilde{\mathbf{A}} = [\tilde{\mathbf{A}}_{Pr}, \tilde{\mathbf{A}}_{Sd}]$  where  $\tilde{\mathbf{A}}_{Pr}$  is an invertible matrix. Let  $\boldsymbol{\mu} := [\mu_i^\circ + \text{RT} \ln \xi_i]_{i=1, \dots, N}$ , since  $\mathbf{S}^T \boldsymbol{\mu} = \mathbf{0}$  and  $\ker \mathbf{S}^T = \text{Im } \mathbf{A}^T$  from (4), there exists  $\mathbf{y} \in \mathbb{R}^M$  such that

$$\boldsymbol{\mu} = \mathbf{A}^T \mathbf{y} \Leftrightarrow \mathbf{P}^T \boldsymbol{\mu} = \tilde{\mathbf{A}}^T \mathbf{y} \Leftrightarrow \tilde{\mathbf{S}}^T \mathbf{P}^T \boldsymbol{\mu} = \mathbf{0},$$

where

$$\tilde{\mathbf{S}} := \begin{bmatrix} \tilde{\mathbf{A}}_{Pr}^{-1} \tilde{\mathbf{A}}_{Sd} \\ -\mathbf{I}_{Sd} \end{bmatrix}. \quad (38)$$

Then one readily checks that  $(\mathbf{n}, \boldsymbol{\mu})$  solves

$$\begin{cases} \mathbf{A}\mathbf{n} = \mathbf{b}, \\ \mathbf{S}^T \boldsymbol{\mu} = \mathbf{0} \end{cases}$$

if and only if  $(\tilde{\mathbf{n}}, \tilde{\boldsymbol{\mu}}) = (\mathbf{P}^T \mathbf{n}, \mathbf{P}^T \boldsymbol{\mu})$  solves

$$\begin{cases} \tilde{\mathbf{A}}\tilde{\mathbf{n}} = \mathbf{b}, \\ \tilde{\mathbf{S}}^T \tilde{\boldsymbol{\mu}} = \mathbf{0}, \end{cases}$$

from which a counterpart to (12) can be derived.  $\square$

**Lemma 3.4.** *Under assumption (H2), the vector  $\boldsymbol{\xi} = (\xi^\alpha)_{\alpha=1, \dots, N_{Ph}}$ , satisfying the system (12), is unique.*

To establish this lemma, we require the following result from information theory.

**Lemma 3.5** (Gibb's inequality). *Suppose that  $\boldsymbol{\xi} = (\xi_i)_{i=1, \dots, n}$  and  $\bar{\boldsymbol{\xi}} = (\bar{\xi}_i)_{i=1, \dots, n}$  are two discrete probability distributions then*

$$\sum_{i=1}^n \xi_i (\ln \xi_i - \ln \bar{\xi}_i) \geq 0,$$

with equality if and only if  $\xi_i = \bar{\xi}_i$ ,  $i = 1, \dots, n$ .

*Proof of Lemma 3.4.* Let  $\mathcal{X} = (\xi^\alpha, s_\alpha, r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  and  $\bar{\mathcal{X}} = (\bar{\xi}^\alpha, \bar{s}_\alpha, \bar{r}_\alpha)_{\alpha=1, \dots, N_{Ph}}$  be two solutions of (12), then

$$\mathbf{A}[s_\alpha \xi^\alpha - \bar{s}_\alpha \bar{\xi}^\alpha]_{\alpha=1, \dots, N_{Ph}} = 0, \quad (39)$$

$$\mathbf{S}^T [\ln \xi^\alpha - \ln \bar{\xi}^\alpha]_{\alpha=1, \dots, N_{Ph}} = 0. \quad (40)$$

Due to (4), we have the relation  $\ker \mathbf{S}^T = (\ker \mathbf{A})^\perp$ , therefore

$$\sum_{\alpha=1}^{N_{Ph}} \langle s_\alpha \xi^\alpha - \bar{s}_\alpha \bar{\xi}^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle = 0. \quad (41)$$

We will prove that the sum (41) is composed of non-negative terms. For each phase  $\alpha$ , there are three cases:

- if  $s_\alpha = \bar{s}_\alpha = 0$ , then  $\langle s_\alpha \xi^\alpha - \bar{s}_\alpha \bar{\xi}^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle = 0$ ;
- if  $s_\alpha = \bar{s}_\alpha > 0$ , then  $\langle s_\alpha \xi^\alpha - \bar{s}_\alpha \bar{\xi}^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle \geq 0$  since the logarithm is an increasing function;
- if  $s_\alpha \neq \bar{s}_\alpha$ , we can assume, without loss of generality, that  $s_\alpha > \bar{s}_\alpha$ . This implies that  $s_\alpha > 0$  and then  $r_\alpha = 0$ . We can rewrite the term in the sum (41) as

$$\bar{s}_\alpha \langle \xi^\alpha - \bar{\xi}^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle + (s_\alpha - \bar{s}_\alpha) \langle \xi^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle =: \underbrace{\bar{s}_\alpha \mathcal{A}_\alpha}_{\geq 0} + \underbrace{(s_\alpha - \bar{s}_\alpha) \mathcal{B}_\alpha}_{> 0}. \quad (42)$$

Applying the convexity inequality:

$$x(\ln x - \ln y) \geq x - y,$$

to  $\mathcal{B}_\alpha$  yields

$$\mathcal{B}_\alpha = \langle \xi^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle \geq \sum_{i \in \sigma^{-1}(\alpha)} \xi_i - \bar{\xi}_i = \bar{r}_\alpha - r_\alpha = \bar{r}_\alpha \geq 0 \quad (43)$$

since  $r_\alpha = 0$ . Consequently, (42) is non-negative.

Therefore, (41) is a vanishing sum of non-negative terms, it follows that for each  $\alpha \in \mathcal{P}$ :

$$\langle s_\alpha \xi^\alpha - \bar{s}_\alpha \bar{\xi}^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle = \bar{s}_\alpha \mathcal{A}_\alpha + (s_\alpha - \bar{s}_\alpha) \mathcal{B}_\alpha = 0. \quad (44)$$

In order to prove the uniqueness of  $\xi^\alpha$ , there are the same three cases to deal with:

- if  $s_\alpha = \bar{s}_\alpha > 0$ , then (44) becomes  $\mathcal{A}_\alpha = 0$ . We introduce the function

$$\mathbf{g}_\alpha(\xi^\alpha) := \sum_{i \in \sigma^{-1}(\alpha)} \xi_i \ln \xi_i,$$

which is strictly convex on  $\mathbb{R}_{>0}^{N_\alpha}$ . We can then rewrite  $\mathcal{A}_\alpha$  as

$$\mathcal{A}_\alpha = \langle \xi^\alpha - \bar{\xi}^\alpha, \nabla \mathbf{g}_\alpha(\xi^\alpha) - \nabla \mathbf{g}_\alpha(\bar{\xi}^\alpha) \rangle.$$

The strict convexity of the function  $\mathbf{g}_\alpha$  on the set  $\mathbb{R}_{>0}^{N_\alpha}$  implies the strict monotonicity of its gradient  $\nabla \mathbf{g}_\alpha$ , consequently

$$\mathcal{A}_\alpha = 0 \quad \Rightarrow \quad \xi^\alpha = \bar{\xi}^\alpha;$$

- if  $s_\alpha > \bar{s}_\alpha$ , then  $s_\alpha > 0$  we have seen that  $r_\alpha = 0$  and

$$\underbrace{\bar{s}_\alpha \mathcal{A}_\alpha}_{\geq 0} + \underbrace{(s_\alpha - \bar{s}_\alpha) \mathcal{B}_\alpha}_{\geq 0} = 0. \quad (45)$$

As a result, both terms in the sum (45) are zero, and therefore,  $\mathcal{B}_\alpha = 0$ . From (43), we deduce that  $\bar{r}_\alpha = 0$ . It follows that  $\langle \xi^\alpha, \mathbf{1} \rangle = 1$  and  $\langle \bar{\xi}^\alpha, \mathbf{1} \rangle = 1$ . This result allows us to apply Gibb's inequality (Lemma 3.5), which states:

$$\mathcal{B}_\alpha = \langle \xi^\alpha, \ln \xi^\alpha - \ln \bar{\xi}^\alpha \rangle = 0 \quad \Leftrightarrow \quad \xi^\alpha = \bar{\xi}^\alpha;$$

- if  $s_\alpha = \bar{s}_\alpha = 0$ , we cannot conclude about the uniqueness of  $\xi^\alpha$  from (44) but Lemma 3.3 guarantees that there are at least  $M$  primary species present in  $\mathcal{X}$  and that the formula matrix  $\mathbf{A}$  can be reformulated with these  $M$  species as the first  $M$  columns. Let  $\xi_{P_r}$  and  $\bar{\xi}_{P_r}$  the corresponding vectors in  $\mathcal{X}$  and  $\bar{\mathcal{X}}$  respectively. We have established in Proposition 3.1 that  $\xi^\alpha > 0$  for each phase  $\alpha$ . This implies that if a species is present in a phase  $\alpha$ , then  $s_\alpha > 0$ , indicating that the entire phase is present. It follows that the phases where  $\xi_{P_r}$  belongs are all present, and from the previous cases we can conclude that  $\xi_{P_r} = \bar{\xi}_{P_r}$ . Let  $\mathbf{d} := \mathbf{S}^T \boldsymbol{\mu}^\circ / (\text{RT})$ , it follows from (12) that for each  $i \in \sigma^{-1}(\alpha)$ ,

$$\begin{aligned} \left( \mathbf{d} + \tilde{\mathbf{A}}_{Sd}^T \tilde{\mathbf{A}}_{P_r}^{-T} [\ln \xi_{P_r}] \right)_i &= \ln \xi_i^\alpha, \\ \left( \mathbf{d} + \tilde{\mathbf{A}}_{Sd}^T \tilde{\mathbf{A}}_{P_r}^{-T} [\ln \bar{\xi}_{P_r}] \right)_i &= \ln \bar{\xi}_i^\alpha, \end{aligned}$$

and since  $\xi_{P_r} = \bar{\xi}_{P_r}$ ,

$$\left( \mathbf{d} + \tilde{\mathbf{A}}_{Sd}^T \tilde{\mathbf{A}}_{P_r}^{-T} [\ln \xi_{P_r}] \right)_i = \ln \bar{\xi}_i^\alpha,$$

which implies that  $\xi^\alpha = \bar{\xi}^\alpha$ .

□

With these two lemmas, we can now demonstrate Proposition 3.2.

*Proof of Proposition 3.2.* In Lemma 3.4, we have proved that under assumption (H2), the vector  $\xi$  is unique. The uniqueness of  $(r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  follows from (12c). To demonstrate the uniqueness of  $(s_\alpha)_{\alpha=1, \dots, N_{Ph}}$ , let  $(\bar{s}_\alpha)_{\alpha=1, \dots, N_{Ph}}$  be another solution. According to equation (12a), we can derive the following relation:

$$\sum_{\alpha \in \Gamma_{\mathcal{X}}} (s_\alpha - \bar{s}_\alpha) \mathbf{A}^\alpha \xi^\alpha = \mathbf{0}.$$

Therefore the uniqueness is satisfied since  $\{\mathbf{A}^\alpha \xi^\alpha\}_{\alpha \in \Gamma_{\mathcal{X}}}$  is linearly independent. □

## 4 Efficient resolution with parametrized Newton’s method

This section focuses on the development of numerical algorithms designed to address the complexities of multiphase chemical equilibrium problems. The core objective is to determine the variables  $(\xi^\alpha, s_\alpha, r_\alpha)_{\alpha=1, \dots, N_{Ph}}$  that satisfy the system of equations (12) using Newton’s method on a new parametrized system which is computationally easier to solve. This system presents two primary challenges:

- **Nonlinearities of the logarithm:** the presence of logarithmic terms introduces stiff nonlinearity, complicating the convergence of iterative methods.
- **Complementarity problem:** the constraints inherent in the system create a complementarity problem (12d)–(12e), necessitating adapted techniques for resolution.

In the first part of this section, we recall the parametrization technique introduced in [18] which was successfully applied to single-phase chemical equilibrium in order to address the issues caused by the logarithmic terms.

In the second part, we propose an innovative approach, referred to as complementarity parametrization, that integrates the complementarity equations directly into the system through parametrization techniques. This approach offers a new framework for solving these problems while maintaining essential conditions along the Newton iterations.

In the third and final part, we establish the invertibility of the Jacobian matrix associated with the combined log and complementarity parametrizations at equilibrium, focusing on a simple two-phase system. This result provides a theoretical foundation for the robustness of the proposed numerical approach.

### 4.1 Addressing log nonlinearity with parametrization

The parametrization technique consists in introducing a fictitious variable  $\tau_i$  for each species and two functions  $X : \mathbb{R} \rightarrow \mathbb{R}$  and  $Y : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$Y(\tau_i) = \ln(X(\tau_i)).$$

These functions allow us to parameterize the graph

$$\{(\xi, y) \in \mathbb{R}^2 \mid y = \ln \xi\},$$

which, up a rescaling and a shift, encodes the relation between the chemical potentials and the (extended) mole fractions. There are various possibilities to define these functions, but the choice made here is to use the *switch* function:

$$(X(\tau), Y(\tau)) = \begin{cases} (\exp(\tau), \tau), & \text{if } \tau < 0, \\ (\tau + 1, \ln(\tau + 1)), & \text{if } \tau \geq 0. \end{cases} \quad (46)$$

Applying this parametrization to (12a)–(12c) gives:

$$\sum_{\alpha=1}^{N_{Ph}} s_\alpha \mathbf{A}^\alpha \mathbf{X}(\boldsymbol{\tau}^\alpha) - b = 0, \quad (47a)$$

$$\mathbf{S}^T [\mu_i^\circ / (RT) + Y(\tau_i)]_{i=1, \dots, N} = \mathbf{0}, \quad (47b)$$

$$\langle \mathbf{X}(\boldsymbol{\tau}^\alpha), \mathbf{1} \rangle + r_\alpha - 1 = 0, \quad (\alpha = 1, \dots, N_{Ph}). \quad (47c)$$

The associated block in the Jacobian is written as:

$$\begin{bmatrix} s_1 \mathbf{A}^1 \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^1)\} & \dots & s_{N_{Ph}} \mathbf{A}^{N_{Ph}} \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^{N_{Ph}})\} & \mathbf{A}\mathbf{X}(\boldsymbol{\tau}) & \mathbf{0} \\ & & \mathbf{S}^T \text{diag}\{Y'(\boldsymbol{\tau})\} & \mathbf{0} & \mathbf{0} \\ & & X'(\boldsymbol{\tau}^1)^T & & \\ & & & \ddots & \\ & & & & X'(\boldsymbol{\tau}^{N_{Ph}})^T & \mathbf{0} & \mathbf{I}_{N_{Ph}} \end{bmatrix}.$$

More details about applying parametrization techniques to chemical equilibria can be found in [18].

## 4.2 Solving the complementarity problem with parametrization

The complementarity problem associated with system (12) is defined by the following conditions:

$$s_\alpha r_\alpha = 0 \quad \text{and} \quad s_\alpha, r_\alpha \geq 0. \quad (48)$$

This formulation indicates that for each pair of non-negative variables  $s_\alpha$  and  $r_\alpha$ , at least one of them must be zero, which is a characteristic of complementarity conditions. Complementarity problems encompass a broad class of mathematical optimization challenges and have generated extensive literature, as noted in works such as those by Acary and Brogliato [1] and Facchinei and Pang [12, 13]. Graphically, the complementarity problem can be visualized as the intersection of the non-negative axes in a Cartesian coordinate system, specifically represented by the two demi-axes  $\{s_\alpha \in \mathbb{R} \mid s_\alpha \geq 0\}$  and  $\{r_\alpha \in \mathbb{R} \mid r_\alpha \geq 0\}$ .

The resolution of the complementarity problem using Newton's method presents two significant challenges:

- Non-differentiability at the origin: the case where  $s_\alpha = r_\alpha$  is not differentiable, which complicates the definition of the Jacobian for the system.
- Nonnegativity constraint violation: the nonnegativity constraint is not inherently preserved by the Newton iterates, which can lead to solutions that fall outside the feasible region.

In this paper, we propose a novel approach to the complementarity problem based on parametrization techniques. In this method, referred to as complementarity parametrization method, the complementarity equations are integrated into the system, and the relationships between the complementarity variables are maintained through a parameter and two parametric functions. This approach provides a new framework for solving complementarity problems while ensuring that the essential conditions are met throughout the iterative process.

The complementarity problem (48) can be described using the min function as

$$\min(s_\alpha, r_\alpha) = 0.$$

Similarly to our approach of parametrization for the relationship  $y = \ln \xi$ , the graph

$$\mathcal{T} = \{(s, r) \in \mathbb{R}^2 \mid \min(s, r) = 0\}$$

is parameterized by two monotonic piecewise continuously differentiable functions,  $S : \mathbb{R} \rightarrow \mathbb{R}$  and  $R : \mathbb{R} \rightarrow \mathbb{R}$ , such that

$$s_\alpha = S(\eta_\alpha) \quad \text{and} \quad r_\alpha = R(\eta_\alpha).$$

These functions are defined to satisfy:

$$\min(S(\eta), R(\eta)) = 0,$$

which implies that  $\mathcal{T} = (S, R)(\mathbb{R})$ . To ensure proper parametrization and avoid singularities, the functions  $S$  and  $R$  must satisfy the following conditions, there exists  $M > 0$ ,  $\varepsilon > 0$  such that for each  $\eta \in \mathbb{R}$ :

(PC1)  $S(\eta)R(\eta) = 0$  and  $S(\eta), R(\eta) \geq 0$  with  $S$  and  $R$  monotonic and  $M$ -Lipschitz continuous;

(PC2) the derivatives  $S'$  and  $R'$  are bounded and piecewise continuous with  $0 \leq S'(\eta) \leq M$  and  $-M \leq R'(\eta) \leq 0$ ;

$$(PC3) \quad S'(\eta) - R'(\eta) > \varepsilon.$$

We define a parametrization as admissible if it fulfills conditions (PC1)–(PC3). To define these specific functions, we require that  $S$  and  $R$  are both monotonic and  $M$ -Lipschitz continuous, which implies that they are differentiable almost everywhere, as established by Rademacher's theorem. For the points at which  $S$  and  $R$  are not differentiable, we assign specific values to  $S'$  and  $R'$ . In practice,  $S'$  and  $R'$  exhibit piecewise continuity. Furthermore, we ensure that  $|S'|$  and  $|R'|$  are upper semi-continuous by choosing the largest possible value (in magnitude) for both  $S'$  and  $R'$ . We also impose a normalization condition:

$$\max(|S'(\eta)|, |R'(\eta)|) = 1, \quad \forall \eta \in \mathbb{R}.$$

Condition (PC1)–(PC3) implies the existence of a point  $\eta^0 \in \mathbb{R}$  such that  $S(\eta^0) = R(\eta^0) = 0$ . By selecting  $\eta^0 = 0$  and enforcing  $S' \geq 0$ ,  $R' \leq 0$ , we derive the following implications:

$$\begin{aligned} \eta > 0 &\Rightarrow S(\eta) > 0 \text{ and } R(\eta) = 0 \Rightarrow R'(\eta) = 0, \\ \eta < 0 &\Rightarrow S(\eta) = 0 \text{ and } R(\eta) > 0 \Rightarrow S'(\eta) = 0. \end{aligned}$$

From the normalization condition, we have:

$$\max(|S'(\eta)|, |R'(\eta)|) = 1 \Leftrightarrow \max(S'(\eta), -R'(\eta)) = 1$$

This leads to the following conclusions:

$$\begin{aligned} \eta > 0 &\Rightarrow S'(\eta) = 1 \Rightarrow S(\eta) = \eta, \\ \eta < 0 &\Rightarrow R'(\eta) = -1 \Rightarrow R(\eta) = -\eta. \end{aligned}$$

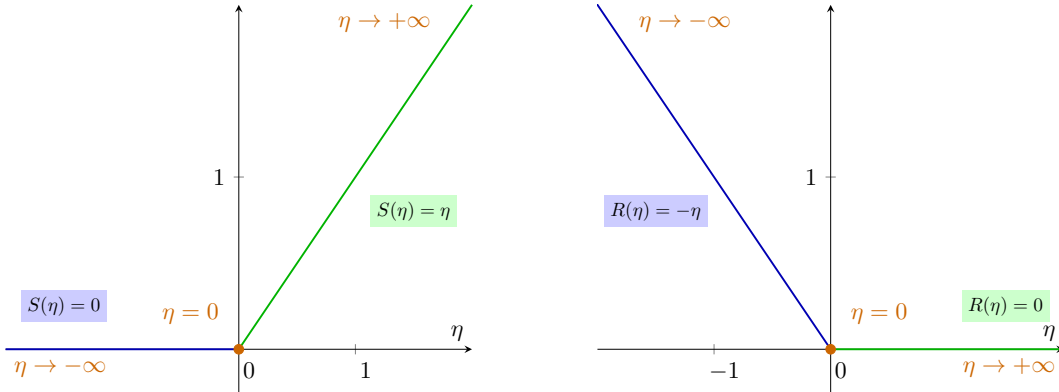
Thus, we define the complementarity parametrization method as follows:

$$S(\eta) = \begin{cases} 0 & \text{if } \eta < 0, \\ \eta & \text{if } \eta \geq 0, \end{cases} \quad \text{and} \quad R(\eta) = \begin{cases} -\eta & \text{if } \eta \leq 0, \\ 0 & \text{if } \eta > 0. \end{cases}$$

These functions are illustrated in Figure 1. Furthermore, the derivatives of these functions are defined as:

$$S'(\eta) = \begin{cases} 0 & \text{if } \eta < 0, \\ 1 & \text{if } \eta \geq 0, \end{cases} \quad \text{and} \quad R'(\eta) = \begin{cases} -1 & \text{if } \eta \leq 0, \\ 0 & \text{if } \eta > 0. \end{cases}$$

At the point  $\eta = 0$ , we have chosen to set both  $|S'(0)|$  and  $|R'(0)|$  equal to 1 to ensure the invertibility of the Jacobian matrix (see Proposition 4.1).



**Figure 1:** Functions  $S$  and  $R$  for the complementarity parametrization method.

These functions lead to a reformulation of system (12) which becomes: find  $(\xi^\alpha, \eta_\alpha)_{\alpha=1, \dots, N_{Ph}}$  such that:

$$\begin{aligned} \sum_{\alpha=1}^{N_{Ph}} S(\eta_\alpha) \mathbf{A}^\alpha \xi^\alpha - b &= 0, \\ \mathbf{S}^T [\mu_i^\circ / (\text{RT}) + \ln \xi_i]_{i=1, \dots, N} &= \mathbf{0}, \\ \langle \xi^\alpha, \mathbf{1} \rangle + R(\eta_\alpha) - 1 &= 0, \quad (\alpha = 1, \dots, N_{Ph}). \end{aligned}$$

Using the parametrization of the log, the system to solve is: find  $(\boldsymbol{\tau}^\alpha, \eta_\alpha)_{\alpha=1, \dots, N_{Ph}}$  such that:

$$\begin{aligned} \sum_{\alpha=1}^{N_{Ph}} S(\eta_\alpha) \mathbf{A}^\alpha \mathbf{X}(\boldsymbol{\tau}^\alpha) - b &= 0, \\ \mathbf{S}^T [\mu_i^\circ / (RT) + Y(\tau_i)]_{i=1, \dots, N} &= \mathbf{0}, \\ \langle \mathbf{X}(\boldsymbol{\tau}^\alpha), \mathbf{1} \rangle + R(\eta_\alpha) - 1 &= 0, \quad (\alpha = 1, \dots, N_{Ph}). \end{aligned}$$

The resulting Jacobian is given by:

$$\begin{bmatrix} S(\eta_1) \mathbf{A}^1 \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^1)\} & \dots & S(\eta_{N_{Ph}}) \mathbf{A}^{N_{Ph}} \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^{N_{Ph}})\} & S'(\eta_1) \mathbf{A}^1 \mathbf{X}(\boldsymbol{\tau}^1) & \dots & S'(\eta_{N_{Ph}}) \mathbf{A}^{N_{Ph}} \mathbf{X}(\boldsymbol{\tau}^{N_{Ph}}) \\ & & \mathbf{S}^T \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau})\} & & & \mathbf{0} \\ & & \mathbf{X}'(\boldsymbol{\tau}^1)^T & & & R'(\eta_1) \\ & & & \ddots & & \\ & & & & \mathbf{X}'(\boldsymbol{\tau}^{N_{Ph}})^T & R'(\eta_{N_{Ph}}) \end{bmatrix}.$$

### 4.3 Invertibility of the complementarity parametrization system

We consider the system combining the logarithmic parametrization with the complementarity parametrization method. To study the invertibility of the associated Jacobian, we focus on the simpler two-phase case: find  $\boldsymbol{\tau} = (\boldsymbol{\tau}^\alpha, \boldsymbol{\tau}^\beta)$  and  $\boldsymbol{\eta} = (\eta_\alpha, \eta_\beta)$  such that:

$$\begin{aligned} \mathbf{A} \begin{bmatrix} S(\eta_\alpha) \mathbf{X}(\boldsymbol{\tau}^\alpha) \\ S(\eta_\beta) \mathbf{X}(\boldsymbol{\tau}^\beta) \end{bmatrix} - \mathbf{b} &= \mathbf{0} \\ \mathbf{S}^T \begin{bmatrix} \mathbf{Y}(\boldsymbol{\tau}^\alpha) \\ \mathbf{Y}(\boldsymbol{\tau}^\beta) \end{bmatrix} + \mathbf{S}^T \boldsymbol{\mu}^\circ / (RT) &= \mathbf{0} \\ \mathbf{1}^T \mathbf{X}(\boldsymbol{\tau}^\alpha) + R(\eta_\alpha) - 1 &= 0 \\ \mathbf{1}^T \mathbf{X}(\boldsymbol{\tau}^\beta) + R(\eta_\beta) - 1 &= 0 \end{aligned} \quad (49)$$

We now prove the invertibility of the Jacobian  $\mathbf{J}(\boldsymbol{\tau}, \boldsymbol{\eta})$  under the assumptions stated earlier.

**Proposition 4.1.** *Assuming (H2) holds and the uniqueness condition from Proposition 3.2 is satisfied, let  $(\boldsymbol{\tau}, \boldsymbol{\eta})$  denote the unique solution to (49). Then, the Jacobian matrix  $\mathbf{J}(\boldsymbol{\tau}, \boldsymbol{\eta})$  is nonsingular.*

*Proof.* Let  $\delta\boldsymbol{\tau} = (\delta\boldsymbol{\tau}^\alpha, \delta\boldsymbol{\tau}^\beta) \in \mathbb{R}^N$  and  $\delta\boldsymbol{\eta} = (\delta\eta_\alpha, \delta\eta_\beta) \in \mathbb{R}^2$  such that  $\mathbf{J}(\boldsymbol{\tau}, \boldsymbol{\eta})(\delta\boldsymbol{\tau}, \delta\boldsymbol{\eta})^T = \mathbf{0}$ , then:

$$\mathbf{A} \begin{bmatrix} S(\eta_\alpha) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\alpha)\} \delta\boldsymbol{\tau}^\alpha + S'(\eta_\alpha) \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta\eta_\alpha \\ S(\eta_\beta) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\beta)\} \delta\boldsymbol{\tau}^\beta + S'(\eta_\beta) \mathbf{X}(\boldsymbol{\tau}^\beta) \delta\eta_\beta \end{bmatrix} = \mathbf{0}, \quad (50a)$$

$$\mathbf{S}^T \begin{bmatrix} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\alpha)\} \delta\boldsymbol{\tau}^\alpha \\ \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \delta\boldsymbol{\tau}^\beta \end{bmatrix} = \mathbf{0}, \quad (50b)$$

$$\mathbf{X}'(\boldsymbol{\tau}^\alpha)^T \delta\boldsymbol{\tau}^\alpha + R'(\eta_\alpha) \delta\eta_\alpha = 0, \quad (50c)$$

$$\mathbf{X}'(\boldsymbol{\tau}^\beta)^T \delta\boldsymbol{\tau}^\beta + R'(\eta_\beta) \delta\eta_\beta = 0. \quad (50d)$$

From (4), we have  $\ker \mathbf{S}^T = (\ker \mathbf{A})^\perp$ . Taking the scalar product of the vector in (50a) with the vector in (50b), we get:

$$\begin{aligned} &\left\langle \begin{pmatrix} S(\eta_\alpha) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\alpha)\} \delta\boldsymbol{\tau}^\alpha \\ S(\eta_\beta) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\beta)\} \delta\boldsymbol{\tau}^\beta \end{pmatrix}, \begin{pmatrix} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\alpha)\} \delta\boldsymbol{\tau}^\alpha \\ \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \delta\boldsymbol{\tau}^\beta \end{pmatrix} \right\rangle \\ &= - \left\langle \begin{pmatrix} S'(\eta_\alpha) \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta\eta_\alpha \\ S'(\eta_\beta) \mathbf{X}(\boldsymbol{\tau}^\beta) \delta\eta_\beta \end{pmatrix}, \begin{pmatrix} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\alpha)\} \delta\boldsymbol{\tau}^\alpha \\ \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \delta\boldsymbol{\tau}^\beta \end{pmatrix} \right\rangle. \end{aligned} \quad (51)$$

Using the identity  $\mathbf{X}'(\boldsymbol{\tau}) = \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau})\} \mathbf{X}(\boldsymbol{\tau})$ , the right-hand side of (51) becomes:

$$- \left\langle \begin{pmatrix} S'(\eta_\alpha) \mathbf{X}'(\boldsymbol{\tau}^\alpha) \delta\eta_\alpha \\ S'(\eta_\beta) \mathbf{X}'(\boldsymbol{\tau}^\beta) \delta\eta_\beta \end{pmatrix}, \delta\boldsymbol{\tau} \right\rangle. \quad (52)$$

We now distinguish three cases:

**Case 1: Both  $\alpha$  and  $\beta$  are present.** Here,  $\eta_\alpha, \eta_\beta > 0$  and we have:

$$S(\eta_\alpha), S(\eta_\beta) > 0, \quad S'(\eta_\alpha) = S'(\eta_\beta) = 1, \quad R(\eta_\alpha) = R(\eta_\beta) = 0, \quad R'(\eta_\alpha) = R'(\eta_\beta) = 0.$$

From (50c) and (50d), we get:

$$\mathbf{X}'(\boldsymbol{\tau}^\alpha)^T \delta \boldsymbol{\tau}^\alpha = 0 \quad \text{and} \quad \mathbf{X}'(\boldsymbol{\tau}^\beta)^T \delta \boldsymbol{\tau}^\beta = 0.$$

Plugging into (52) and then (51), we obtain:

$$\left\langle \underbrace{\begin{pmatrix} S(\eta_\alpha) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\alpha)\} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\alpha)\} \delta \boldsymbol{\tau}^\alpha & \mathbf{0} \\ \mathbf{0} & S(\eta_\beta) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\beta)\} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \end{pmatrix}}_{:=\mathbf{D}} \delta \boldsymbol{\tau}, \delta \boldsymbol{\tau} \right\rangle = 0,$$

where  $\mathbf{D}$  is positive definite since  $\mathbf{X}'(\boldsymbol{\tau}), \mathbf{Y}'(\boldsymbol{\tau}), \mathbf{S}(\boldsymbol{\eta}) > 0$ . Thus,  $\delta \boldsymbol{\tau} = \mathbf{0}$ . Equation (50a) then gives:

$$\mathbf{A} \begin{bmatrix} \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta \eta_\alpha \\ \mathbf{X}(\boldsymbol{\tau}^\beta) \delta \eta_\beta \end{bmatrix} = \mathbf{0} \quad \Leftrightarrow \quad \mathbf{A}^\alpha \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta \eta_\alpha + \mathbf{A}^\beta \mathbf{X}(\boldsymbol{\tau}^\beta) \delta \eta_\beta = \mathbf{0}.$$

According to Proposition 3.2, the vectors  $\{\mathbf{A}^\alpha \mathbf{X}(\boldsymbol{\tau}^\alpha), \mathbf{A}^\beta \mathbf{X}(\boldsymbol{\tau}^\beta)\}$  form a linearly independent set. Consequently, we must have  $\delta \eta_\alpha = \delta \eta_\beta = 0$ .

**Case 2: Phase  $\alpha$  is present,  $\beta$  is absent with  $R(\eta_\beta) > 0$ .** We have  $\eta_\alpha > 0, \eta_\beta < 0$ , then

$$S(\eta_\alpha) > 0, S(\eta_\beta) = 0, \quad S'(\eta_\alpha) = 1, S'(\eta_\beta) = 0, \quad R(\eta_\alpha) = 0, R(\eta_\beta) > 0, \quad R'(\eta_\alpha) = 0, R'(\eta_\beta) = -1.$$

From (50c) we again have:  $\mathbf{X}'(\boldsymbol{\tau}^\alpha)^T \delta \boldsymbol{\tau}^\alpha = 0$  and (52) vanishes. Then (51) yields:

$$\left\langle \underbrace{S(\eta_\alpha) \text{diag}\{\mathbf{X}'(\boldsymbol{\tau}^\alpha)\} \text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\alpha)\} \delta \boldsymbol{\tau}^\alpha, \delta \boldsymbol{\tau}^\alpha}_{:=\mathbf{D}_\alpha} \right\rangle = 0,$$

where  $\mathbf{D}_\alpha$  is positive definite since  $\mathbf{X}'(\boldsymbol{\tau}^\alpha), \mathbf{Y}'(\boldsymbol{\tau}^\alpha), \mathbf{S}(\eta_\alpha) > 0$ . It follows that  $\delta \boldsymbol{\tau}^\alpha = \mathbf{0}$ . Under assumption (H2), the set of vectors  $\mathbf{A}_\alpha = [\mathbf{a}_i]_{i \in \sigma^{-1}(\alpha)}$  is a spanning set in  $\mathbb{R}^M$ . In particular,  $\mathbf{A}_{Pr}$  is included into  $\mathbf{A}_\alpha$  and from the structure

$$\mathbf{S}^T = [(\mathbf{A}_{Pr}^{-1} \mathbf{A}_{Sd})^T, -\mathbf{I}_{Sd}],$$

it follows that (50b) becomes:

$$-\text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \delta \boldsymbol{\tau}^\beta = \mathbf{0}.$$

Since  $\mathbf{Y}'(\boldsymbol{\tau}^\beta) > 0$ , we have  $\delta \boldsymbol{\tau}^\beta = \mathbf{0}$ . Then, from (50d), it follows that  $\delta \eta_\beta = 0$ . Finally, equation (50a) reduces to:

$$\mathbf{A} \begin{bmatrix} \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta \eta_\alpha \\ \mathbf{0} \end{bmatrix} = \mathbf{0},$$

implying  $\delta \eta_\alpha = 0$  by Proposition 3.2.

**Case 3: Phase  $\alpha$  is present,  $\beta$  is absent with  $R(\eta_\beta) = 0$ .** Here,  $\eta_\alpha > 0, \eta_\beta = 0$  and

$$S(\eta_\alpha) > 0, S(\eta_\beta) = 0, \quad S'(\eta_\alpha) = 1, S'(\eta_\beta) = 0, \quad R(\eta_\alpha) = 0, R(\eta_\beta) = 0, \quad R'(\eta_\alpha) = 0, R'(\eta_\beta) = -1.$$

Equations (50c) and (50d) become:

$$\mathbf{X}'(\boldsymbol{\tau}^\alpha)^T \delta \boldsymbol{\tau}^\alpha = 0 \quad \text{and} \quad \mathbf{X}'(\boldsymbol{\tau}^\beta)^T \delta \boldsymbol{\tau}^\beta = \delta \eta_\beta.$$

Then from (51) and (52), we obtain:

$$\langle \mathbf{D}_\alpha \delta \boldsymbol{\tau}^\alpha, \delta \boldsymbol{\tau}^\alpha \rangle = -(\delta \eta_\beta)^2,$$

which implies  $\delta \boldsymbol{\tau}^\alpha = \mathbf{0}$ , and  $\delta \eta_\beta = 0$  since  $\mathbf{D}_\alpha$  is positive definite. From (50a), we get:

$$\mathbf{A} \begin{bmatrix} \mathbf{X}(\boldsymbol{\tau}^\alpha) \delta \eta_\alpha \\ \mathbf{0} \end{bmatrix} = \mathbf{0} \quad \Rightarrow \quad \delta \eta_\alpha = 0$$

thanks to Proposition 3.2. Reasoning analogously to Case 2, we have:

$$-\text{diag}\{\mathbf{Y}'(\boldsymbol{\tau}^\beta)\} \delta \boldsymbol{\tau}^\beta = \mathbf{0} \quad \Rightarrow \quad \delta \boldsymbol{\tau}^\beta = \mathbf{0}.$$

In all cases, the only solution is the trivial one. Hence, the Jacobian is nonsingular.  $\square$

## 5 Numerical experiments

The following section presents a comprehensive study of numerical experiments designed to evaluate computational methods for solving chemical equilibrium problems in multiphase systems. We examine three test cases of increasing complexity: (1) the simple silica ( $\text{SiO}_2$ ) system already introduced as an illustration in Section 2.1, (2) a carbon dioxide ( $\text{CO}_2$ ) system involving aqueous, mineral, and gaseous phases, and (3) a complex seawater system with numerous species and minerals.

The primary objective of this study is to compare our proposed methods—the parametrized approach for the log formulation and the complementarity parametrization strategy—with established methods from the literature, namely the log-trick for the log formulation and the minimum function (Newton-min), Fischer-Burmeister function (FB), and Interior Point Method (IPM) for the complementarity problem. Details of the literature methods are provided in Appendix C.

Our analysis focuses on the performance and robustness of these methods for chemical equilibrium calculations. We pay particular attention to convergence behavior, iteration counts, and the challenges posed by complementarity conditions in equilibrium modeling. The results provide critical insights into the efficiency and reliability of these methods across diverse chemical scenarios.

### Numerical parameters

The functions  $X$  and  $Y$  used in the complementarity parametrization are those defined by the switching mechanism in Eq. (46). For the IPM, we use the following parameter values:  $\kappa_\nu = 0.5$ ,  $\theta_\nu = 2$ , and  $\gamma_{\min} = 0.1$ .

Let  $\mathcal{F}(\boldsymbol{\mathcal{X}})$  denote the residual function whose root we seek. The convergence criterion for Newton’s method is defined as:

$$\|\mathcal{F}(\boldsymbol{\mathcal{X}}^{(k+1)})\|_\infty \leq 10^{-10} \quad \text{and} \quad \|\boldsymbol{\mathcal{X}}^{(k+1)} - \boldsymbol{\mathcal{X}}^{(k)}\|_\infty \leq 10^{-10},$$

where  $k + 1$  denotes the current Newton iteration.

In the IPM framework, the complementarity variables  $s_\alpha$  and  $r_\alpha$  remain strictly positive and never reach zero. This poses challenges when modeling equilibrium conditions that require both variables to be simultaneously zero. To address this, we apply a stricter convergence criterion to the complementarity residual:

$$\|\mathcal{F}(\boldsymbol{\mathcal{X}}^{(k)})\|_\infty \leq 10^{-10} \quad \text{and} \quad \|\mathcal{F}_{\text{compl}}(\boldsymbol{\mathcal{X}}^{(k)})\|_\infty \leq \epsilon_{64}^2,$$

where  $\epsilon_{64} \approx 2.22 \times 10^{-16}$  is the machine epsilon in double-precision floating-point arithmetic.

### Initialization of Newton’s Method.

We initialize the Newton solver using the following procedure. The initial vectors of quantities are defined as:

$$n_{\text{aq}} = (\mathbf{b}_O, 1, \dots, 1), \quad n_{\text{min}} = 1 \text{ for all minerals}, \quad n_{\text{gas}} = (1, \dots, 1),$$

where  $\mathbf{b}_O$  is the fixed oxygen content, serving as an approximation for the aqueous  $\text{H}_2\text{O}$  quantity. The total amount in each phase is computed as:

$$s_{\text{aq}} = \langle n_{\text{aq}}, \mathbf{1} \rangle, \quad s_{\text{min}} = n_{\text{min}} \text{ for all minerals}, \quad s_{\text{gas}} = \langle n_{\text{gas}}, \mathbf{1} \rangle.$$

This leads to the initial values for the phase composition vector  $\boldsymbol{\xi}$ :

$$\boldsymbol{\xi}_{\text{aq}} = n_{\text{aq}}/s_{\text{aq}}, \quad \boldsymbol{\xi}_{\text{min}} = 1 \text{ for all minerals}, \quad \boldsymbol{\xi}_{\text{gas}} = n_{\text{gas}}/s_{\text{gas}}.$$

The auxiliary variable  $\mathbf{r}$  is initialized uniformly:

$$r_{\text{aq}} = 1, \quad r_{\text{min}} = 1 \text{ for all minerals}, \quad r_{\text{gas}} = 1.$$

For the complementarity parametrization method, the variable  $\boldsymbol{\eta}$  is initialized as:

$$\eta_{\text{aq}} = s_{\text{aq}}, \quad \eta_{\text{min}} = -s_{\text{min}} \text{ for all minerals}, \quad \eta_{\text{gas}} = -s_{\text{gas}}.$$

Finally, the parameter  $\nu$  for IPM is initialized as:

$$\nu = \max \left\{ \min_{\alpha \in \{1, \dots, N_{Ph}\}} (s_\alpha r_\alpha) - 0.1, 0.5 \right\}.$$

### Description of test cases

We consider three chemical equilibrium test cases of increasing complexity:

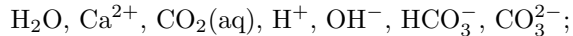
1. SiO<sub>2</sub> system: this simple system consists of four aqueous species:



and one mineral phase: quartz (SiO<sub>2</sub>(quartz)).

2. CO<sub>2</sub> system: this system comprises ten species, divided among three phases:

• Aqueous phase:



• Mineral phase: calcite (CaCO<sub>3</sub>(calcite));

• Gaseous phase: carbon dioxide (CO<sub>2</sub>(g)) and water vapor (H<sub>2</sub>O(g)).

3. Multiphase Seawater: this complex system includes:

• 50 aqueous species;

• 20 minerals, each in its own pure phase;

• 2 gaseous species.

Detailed compositions and thermodynamic data for all three test cases are provided in Appendix A.

### Numerical evaluation

In the two following subsections, we assess the proposed algorithms on the three representative test cases of increasing complexity presented above.

The SiO<sub>2</sub> case, a small-scale chemical system, allows detailed examination of convergence behavior and local phenomena, providing insights into algorithmic mechanisms impractical to obtain in larger systems.

The CO<sub>2</sub> and multiphase Seawater cases follow the same robustness assessment methodology, differing primarily in system complexity. The CO<sub>2</sub> system, relatively simple, serves as a reference for evaluating algorithmic reliability under varied initial guesses and perturbed data. The multiphase Seawater system, highly complex and multiphase, extends this evaluation to large-scale scenarios.

This progressive approach enables both a precise understanding of algorithmic behavior and a rigorous evaluation of practical performance in challenging settings.

### 5.1 SiO<sub>2</sub>: small-scale convergence analysis

For this test case, we have run our algorithms with three different vectors **b** in order to obtain the following subcases:

- **A**: the mineral is present in a large amount;
- **B**: the mineral is absent with condition that are close to precipitation, indicating that both variables of the complementarity conditions vanish at equilibrium;
- **C**: the mineral is absent with conditions that are far from precipitation.

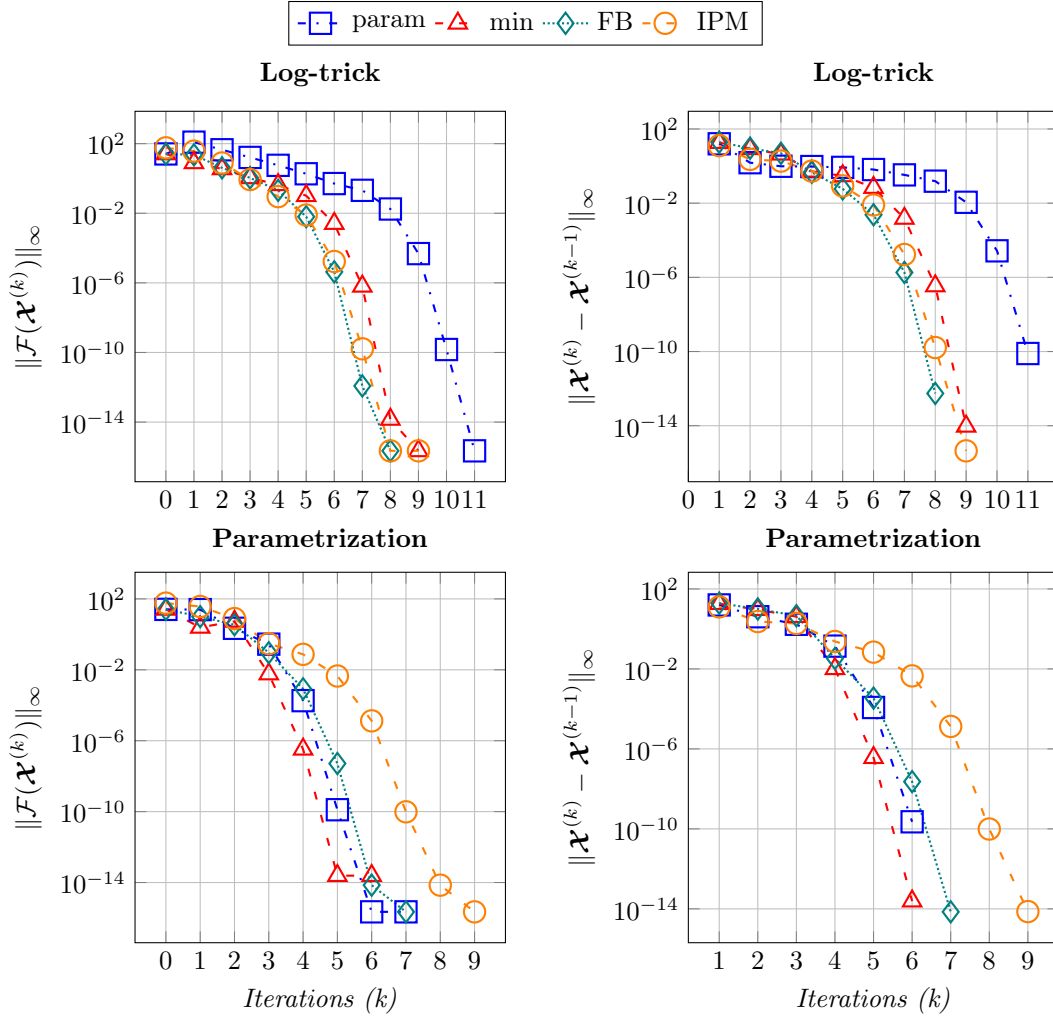
The results obtained from the various combinations of methods are presented in Table 1. The evolution of residuals for setting **A** is illustrated in Figure 2 and the evolution of IPM in setting **B** in Figure 7. Notably, for the methods that converge, local quadratic convergence is observed. These results enable us to draw the following conclusions:

1. In **B**, the IPM is very slow, due to the needed precision on the complementarity equations and to a linear convergence.
2. In **C**, the min function diverges with the parametrization formulation and the FB function is slow with the log-trick and parametrization formulations.
3. The complementarity parametrization method and the IPM approaches are robust. However, the convergence of the IPM is slower.

In the remainder of this section, we will conduct a comprehensive analysis of the diverse behaviors exhibited by algorithms.

Log formulation	Complementarity	Setting		
		A	B	C
Log-trick	param	11	9	13
	min	9	7	10
	FB	8	7	15
	IPM	9	52	12
Param	param	7	6	14
	min	6	6	×
	FB	7	6	14
	IPM	9	52	12

**Table 1:** Number of iterations for Newton’s method.



**Figure 2:** Evolution of residuals for the  $\text{SiO}_2$  text case in configuration **A**.

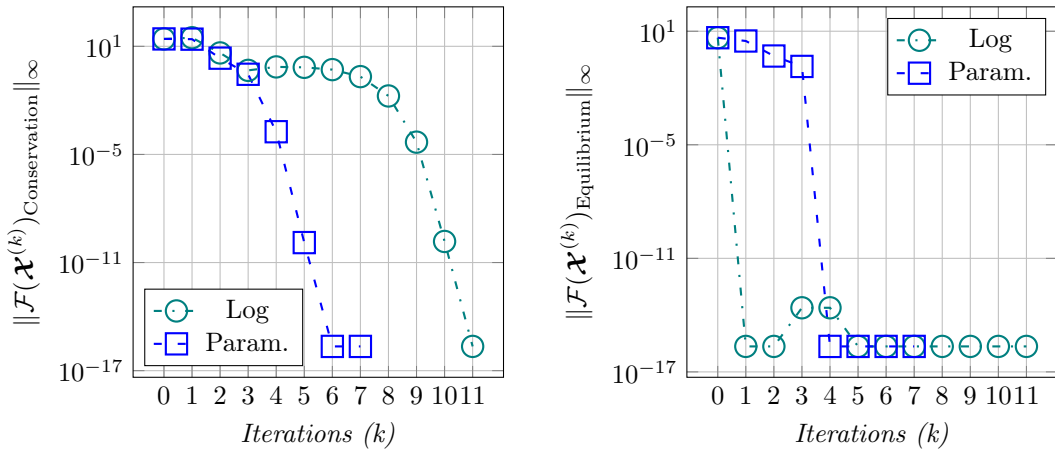
### Comparison between the log-trick and the parametrization with complementarity parametrization

In this part, we compare the iterates produced by the two log-based formulations—the log-trick and the parametrization method—in the context of complementarity-based parametrization. To this end, Figure 3 illustrates the evolution of the norm of the function  $\mathcal{F}$ . The left panel shows the residuals associated with the conservation and mole fraction sum equations, while the right panel focuses on the residuals of the equilibrium equations. Figures 4–5 display the evolution of the iterates for the species  $\text{H}_2\text{O}$  and *Quartz*, along with their mole fractions. Figure 6 presents the corresponding evolution for  $\text{H}^+$ ,  $\text{OH}^-$ , and  $\text{SiO}_2(\text{aq})$ . Table 2 indicates the parametrization function used for  $\text{H}_2\text{O}$  at each iteration.

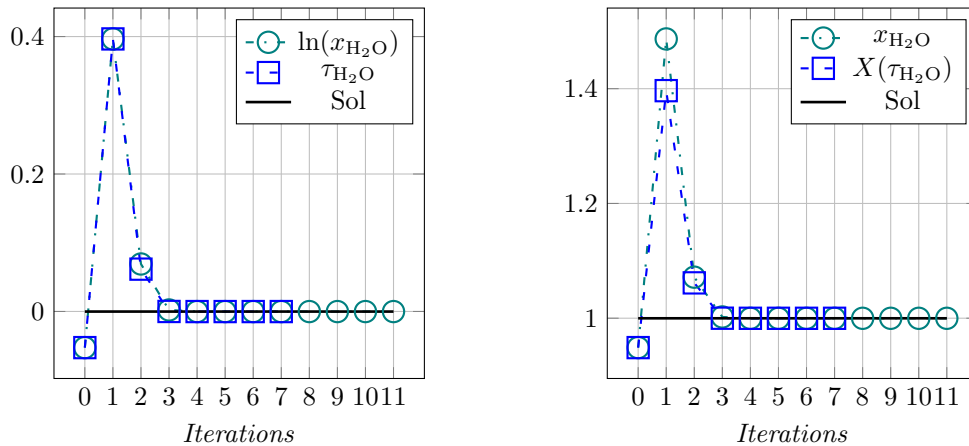
As shown in Figure 3, the decrease in the residual for the parametrization method can be attributed

to the complete resolution of the chemical equilibrium equation from the fourth iteration onward. Specifically, starting at the third iteration, the equilibrium equation becomes linear. It is well known that Newton’s method solves linear equations exactly in a single step. According to Figure 6, the  $\tau_i$  variables—excluding those for  $\text{H}_2\text{O}$  and  $\text{Quartz}$ —are consistently negative, which ensures that the  $Y(\tau_i)$  values remain linear. Consequently, the linearity or non-linearity of the equilibrium equations depends solely on the values of  $\tau_{\text{H}_2\text{O}}$  and  $\tau_{\text{Quartz}}$ . Figure 4 and Table 2 show that  $\tau_{\text{H}_2\text{O}}$  becomes negative starting from the third iteration, and Figure 5 shows that  $\tau_{\text{Quartz}}$  becomes exactly zero at the same point. These conditions together yield a linear equilibrium equation from the third iteration onward.

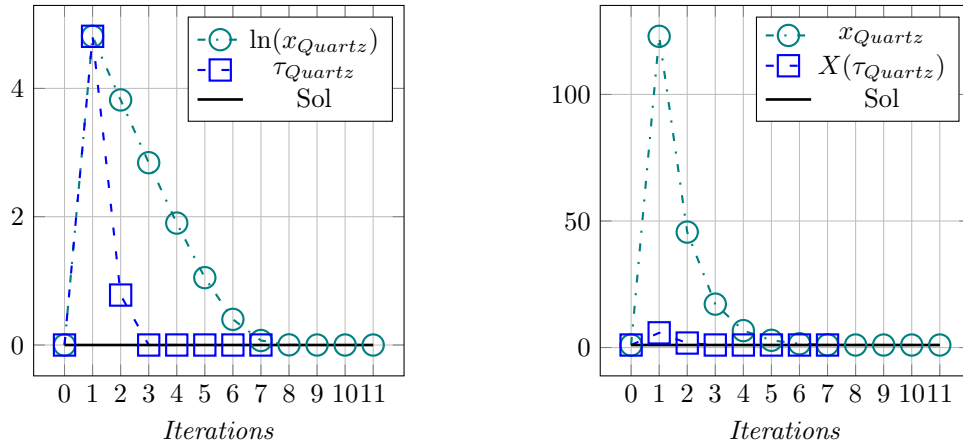
Furthermore, Figure 5 shows that during the first iteration, the value of  $\tau_{\text{Quartz}}$  from the parametrization method coincides with  $\ln(x_{\text{Quartz}})$  obtained using the log-trick approach. However, the corresponding mole fractions,  $X(\tau_{\text{Quartz}})$  and  $x_{\text{Quartz}}$ , differ because  $X(\tau) = \tau + 1$  for positive values of  $\tau$ . This property prevents the mole fraction from becoming excessively large, resulting in a  $\tau_{\text{Quartz}}$  value that is significantly closer to the solution by the second iteration compared to the log-trick. Consequently, in this test case, the parametrization method leads to faster convergence toward the solution.



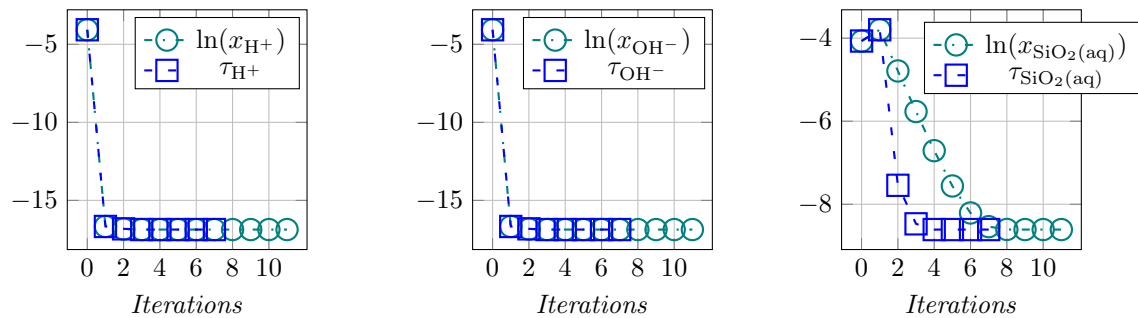
**Figure 3:** Evolution of the residuals for the log-trick and the parametrization methods in the  $\text{SiO}_2$  test case. The left panel shows the residuals of the conservation equations, while the right panel shows those of the equilibrium equations.



**Figure 4:** Evolution of the species  $\text{H}_2\text{O}$  with the log-trick and parametrization methods.



**Figure 5:** Evolution of the *Quartz* with the log-trick and parametrization methods.



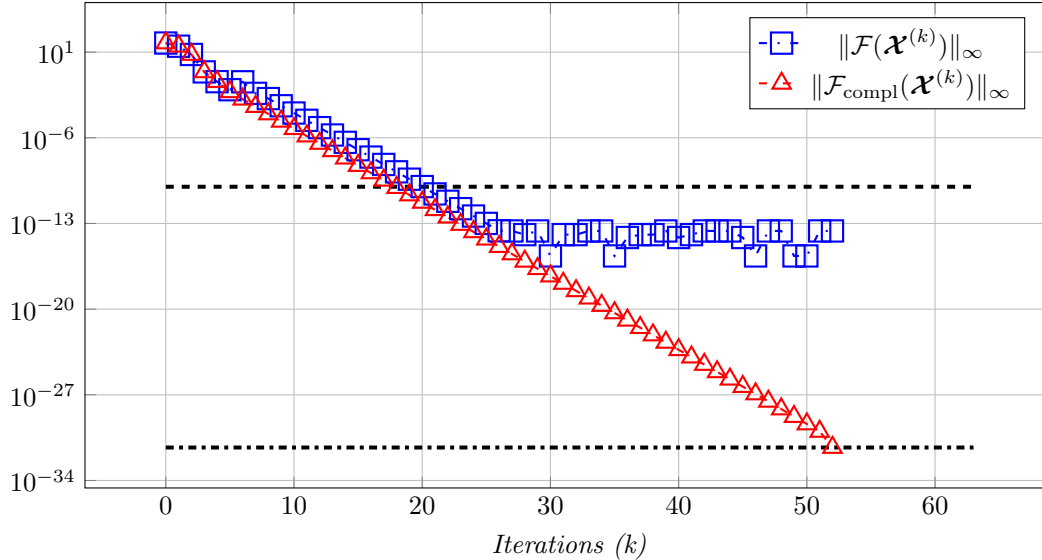
**Figure 6:** Evolution of the species  $\text{H}^+$ ,  $\text{OH}^-$  and  $\text{SiO}_2(\text{aq})$  for the log trick and parametrization method.

Iteration	0	1	2	3	...	7
$\text{sign}(\tau_{\text{H}_2\text{O}})$	-	+	+	-	...	-
$X(\tau_{\text{H}_2\text{O}})$	$\exp \tau_{\text{H}_2\text{O}}$	$\tau_{\text{H}_2\text{O}} + 1$	$\tau_{\text{H}_2\text{O}} + 1$	$\exp \tau_{\text{H}_2\text{O}}$	...	$\exp \tau_{\text{H}_2\text{O}}$
$Y(\tau_{\text{H}_2\text{O}})$	$\tau_{\text{H}_2\text{O}}$	$\ln(\tau_{\text{H}_2\text{O}} + 1)$	$\ln(\tau_{\text{H}_2\text{O}} + 1)$	$\tau_{\text{H}_2\text{O}}$	...	$\tau_{\text{H}_2\text{O}}$

**Table 2:** Evolution of  $X(\tau_{\text{H}_2\text{O}})$  and  $Y(\tau_{\text{H}_2\text{O}})$ .

### IPM slowdown in setting B

Test case **B** is characterized by a uniquely doubly active complementarity constraint associated with the mineral phase. As shown in Figure 7, the residual norm  $\|\mathcal{F}(\boldsymbol{\mathcal{X}}^{(k)})\|_\infty$  reaches the prescribed tolerance at iteration 21. Beyond this point, the residual exhibits a plateau near machine epsilon, persisting until the convergence criterion for the complementarity residual  $\|\mathcal{F}_{\text{compl}}(\boldsymbol{\mathcal{X}}^{(k)})\|_\infty$  is eventually satisfied. Notably, the convergence behavior transitions from quadratic to linear, indicating a degradation in the rate of convergence.

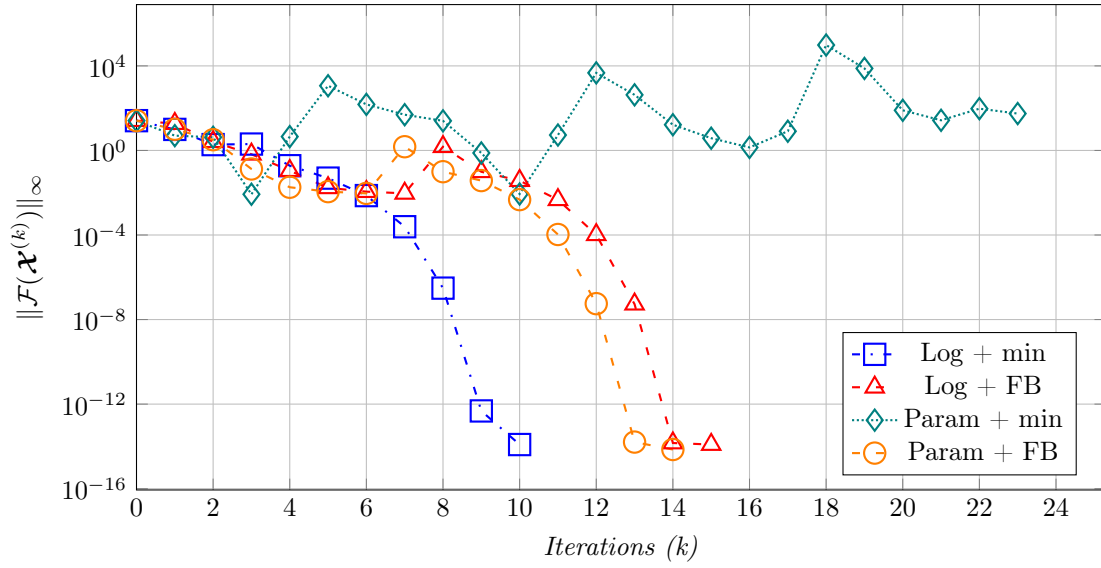


**Figure 7:** Evolution of residuals of IPM for setting **B** with the parametrization formulation.

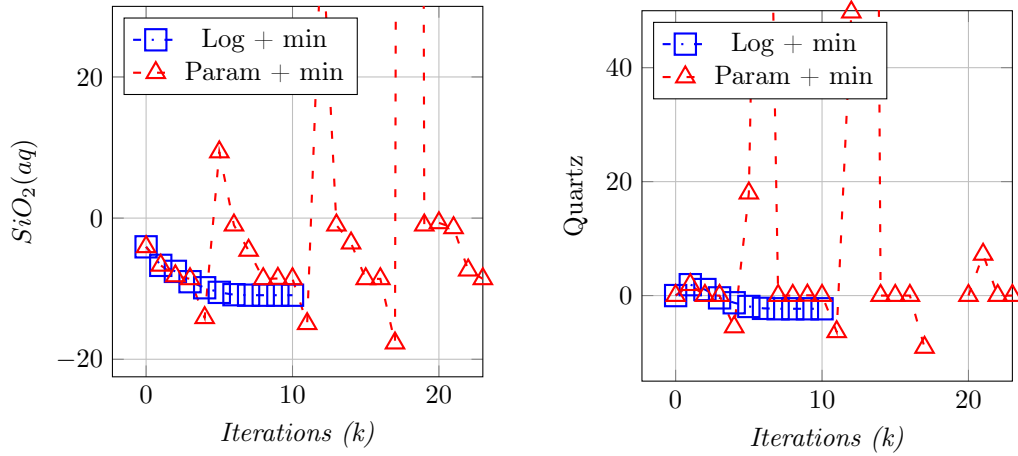
### Slowdown in setting C

The slowdown observed in setting **C**, as shown in Figure 8, can be attributed to a sharp decline in the concentrations of  $\text{SiO}_2$  and Quartz, depicted in Figure 10 in the case of FB function. This abrupt variation fails to effectively reduce the residual, causing the iterates to undershoot the solution before undergoing a large compensatory correction that overshoots it. A comparable behavior occurs with the parametrization based on the min function (Figure 9); however, in this case, the correction fails to recover convergence. Specifically, Newton's method applied to the min-function parametrization terminates prematurely due to a singular Jacobian. At iteration 23, the variable  $n_{\text{aq}}$  is set to zero, leading to a loss of structural information in the Jacobian matrix (53). As a result, the matrix no longer maintains linear independence between rows 3 and 8, and the entries in row 3 approach zero, exacerbating the singularity.

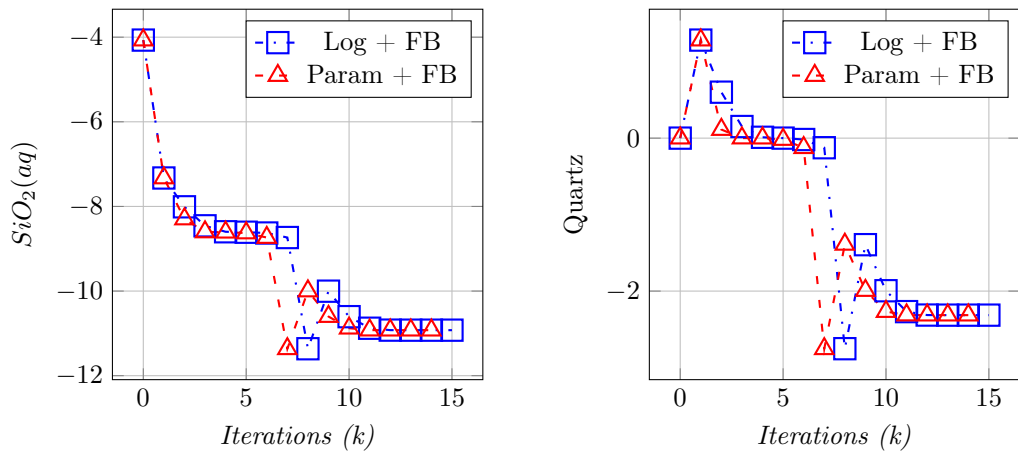
$$\mathcal{J}(\boldsymbol{\mathcal{X}}^{(23)}) \approx \begin{pmatrix} 0.0 & 0.0 & 0.0 & 0.0 & -0.037 & 0.47 & 2.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & -0.018 & 0.00018 & 1.0 & 0.0 & 0.0 \\ \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{3 \times 10^{-22}} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} \\ 1.0 & -1.0 & 0.0 & -1.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & -1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.47 & 4 \times 10^{-8} & 0.00018 & 4 \times 10^{-8} & 0.0 & 0.0 & 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{1.0} & \mathbf{0.0} & \mathbf{0.0} & \mathbf{0.0} \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{pmatrix} \quad (53)$$



**Figure 8:** Evolution of residuals in setting **C**.



**Figure 9:** Evolution of the concentrations of  $\text{SiO}_2(\text{aq})$  and Quartz in setting **C** with the min function.



**Figure 10:** Evolution of the concentrations of  $\text{SiO}_2(\text{aq})$  and Quartz in setting **C** with the FB function.

## 5.2 Robustness analysis across system complexity

In this subsection, we examine the robustness of our proposed methods by analyzing various  $\mathbf{b}$  vectors and different initializations. We will compare the success rates and the number of iterations required for each approach. Through this methodology, we aim to rigorously assess the performance and reliability of our methods under varying conditions.

### Generation of random initialization

To construct a set of random initial guess, we start with a vector  $\mathbf{n}$  such that

$$n_{aq} = (\mathbf{b}_O, 1, \dots, 1), \quad n_{min} = 1, \text{ for all mineral}, \quad n_{gas} = (1, \dots, 1),$$

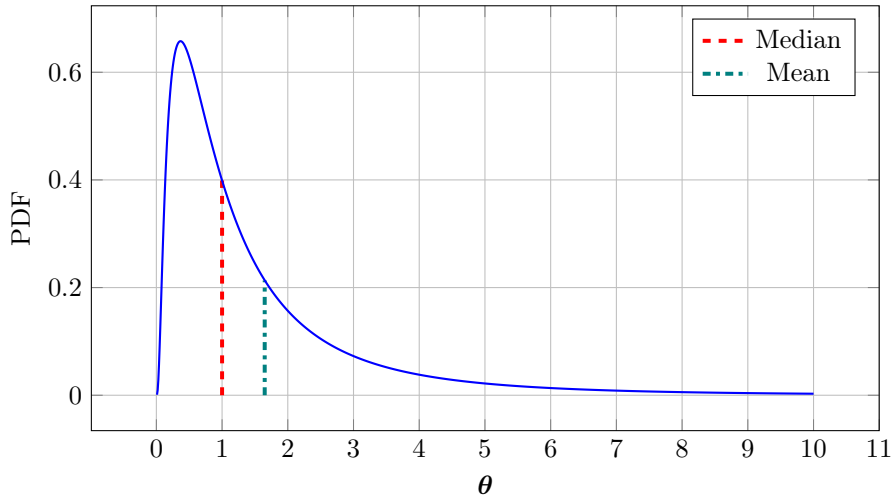
and we generate a perturbed vector  $\tilde{\mathbf{n}}_\theta$  for all species except  $\text{H}_2\text{O}(\text{aq})$  using a log normal distribution characterized by parameters  $(\mu, \sigma)$ :

$$\tilde{\mathbf{n}}_\theta = \mathbf{n} \odot \boldsymbol{\theta} \quad \text{with} \quad \boldsymbol{\theta} \sim \mathcal{LN}(\mu, \sigma^2).$$

where  $\odot$  denotes the Hadamard product. The lognormal distribution ensures that all elements of  $\tilde{\mathbf{n}}_\theta$  remain positive. From the vector  $\mathbf{n}_\theta = [\mathbf{b}_O, \tilde{\mathbf{n}}_\theta]$ , we then define an initial guess  $\boldsymbol{\chi}_\theta^{(0)}$  as detailed in the begin of Section 5. Consequently, we consider a collection of  $K$  random initial guess defined as:

$$\{(\boldsymbol{\chi}_{\theta_1}^{(0)}, \dots, \boldsymbol{\chi}_{\theta_K}^{(0)}) \mid \theta_i \sim \mathcal{LN}(\mu, \sigma^2)\}$$

In our study, we will consider  $K = 1000$  random initial guess characterized by  $\mu = -0.5$  and  $\sigma = 1$  in order to get a median of 1. The corresponding probability density function for  $\boldsymbol{\theta}$  is depicted in Figure 11



**Figure 11:** Probability density function of  $\boldsymbol{\theta} \sim \mathcal{LN}(0, 1)$ .

### Generation of random right-hand side

To construct a set of random vectors for the right-hand side, we start with the solution  $\mathbf{n}^*$  to the minimization problem:

$$\mathbf{n}^* = \arg \min\{G(\mathbf{n}) \mid \mathbf{A}\mathbf{n} = \mathbf{b}\}.$$

To maintain the concept of a diluted solution, we conserve the amount of water from the optimal solution while randomly perturbing the concentrations of the diluted species. We generate a perturbed vector  $\tilde{\mathbf{n}}_\theta$  for all diluted species using a lognormal distribution characterized by parameters  $(\mu, \sigma)$ :

$$\tilde{\mathbf{n}}_\theta = \tilde{\mathbf{n}}^* \odot \boldsymbol{\theta} \quad \text{with} \quad \boldsymbol{\theta} \sim \mathcal{LN}(\mu, \sigma^2).$$

where  $\odot$  denotes the Hadamard product. The lognormal distribution ensures that all elements of  $\tilde{\mathbf{n}}_\theta$  remain positive. From the vector  $\mathbf{n}_\theta = [n_1^*, \tilde{\mathbf{n}}_\theta]$ , we then define a new right-hand side vector as  $\mathbf{b}_\theta = \mathbf{A}\mathbf{n}_\theta$ .

This construction satisfies the non-empty assumption for the set  $\mathcal{M}_{\mathbf{A}, \mathbf{b}_\theta}$ . Consequently, we consider a collection of  $K$  random right-hand sides defined as:

$$\{(\mathbf{b}_{\theta_1}, \dots, \mathbf{b}_{\theta_K}) \mid \theta_i \sim \mathcal{LN}(\mu, \sigma^2)\}$$

### 5.2.1 CO<sub>2</sub>: reference system for robustness evaluation

In this test case, we analyze the sensitivity of the methods to both initialization and variations in the right-hand side vector  $\mathbf{b}$ . Three configurations are considered:

- **A**: only the aqueous phase is present;
- **B**: the aqueous and mineral phases are present;
- **C**: all phases are present.

**Sensitivity to initialization.** The results, shown in Figure 12, include two key metrics: the top panel reports the success rate of each method over 1000 chemical equilibrium computations, while the bottom panel presents the average number of iterations for cases where Newton’s algorithm converges. These results demonstrate the strong robustness of IPM. Notably, the combination of the log-trick with complementarity-based parametrization achieves a comparable success rate, although it requires approximately twice as many iterations to converge compared to the IPM approach.

The relatively poorer performance of the other approaches should be interpreted with care, since the IPM benefits from the line search strategy used in IPOPT, as discussed in Section C.4. It would be worthwhile to test this same line search strategy for the min and FB methods, as they employ the same set of variables as IPM. In contrast, the complementarity parametrization uses only one variable per phase, which introduces different challenges.

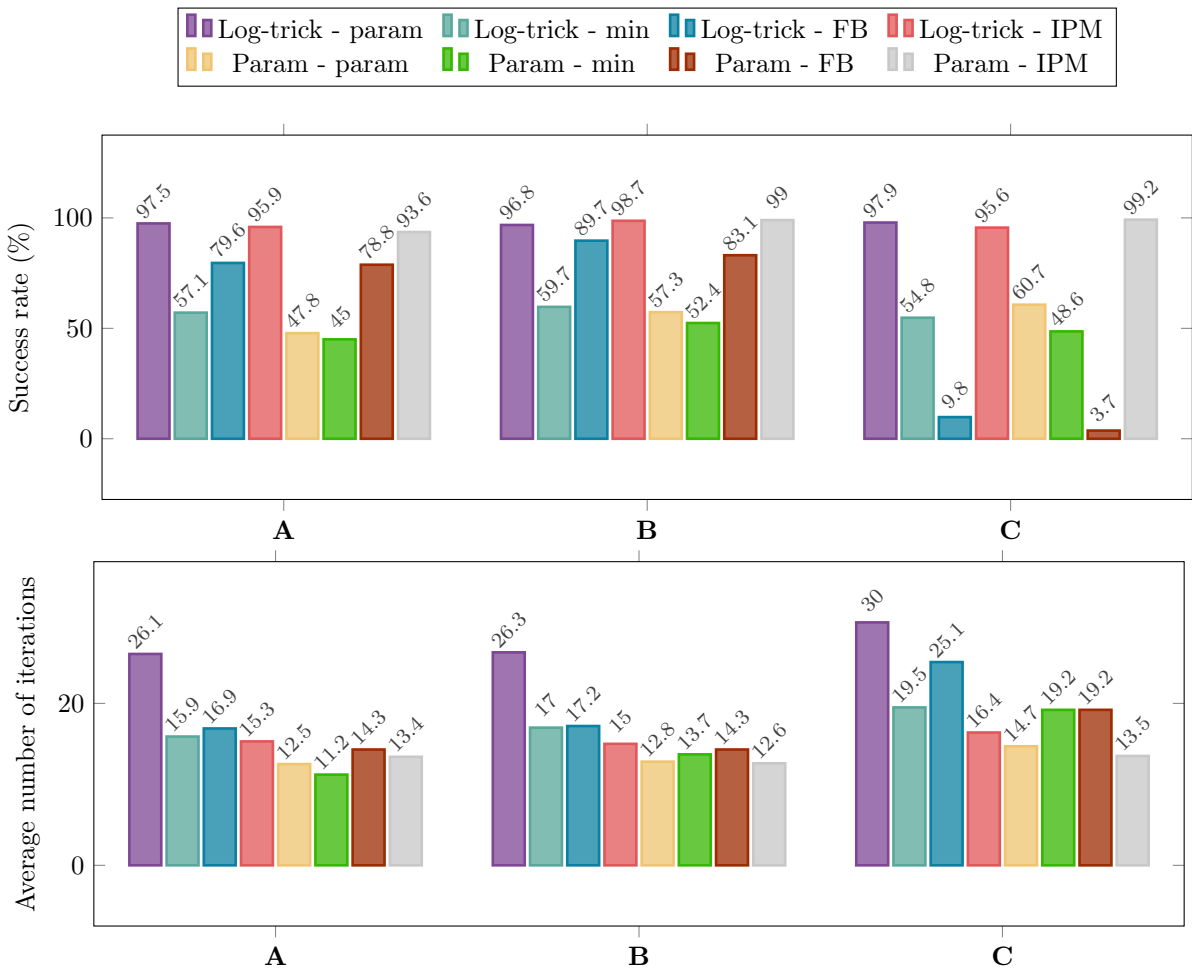
One such challenge is related to invertibility: when a significant portion of the Jacobian becomes null, convergence can be compromised. To address this, we introduce a strategy that "pauses" the update at zero when the sign of a component  $\eta_\alpha$  changes. This approach helps the iterate remains in its current zone at the next iteration, leveraging the particular choice of the derivatives  $S'(0)$  and  $R'(0)$  described in Section 4.2. Although this strategy is applied selectively, we keep using the Newton direction for  $\boldsymbol{\eta}$ . The line search is then modified as follows:

$$\boldsymbol{\eta}^{(k+1)} = \boldsymbol{\eta}^{(k)} + \min_{j \in \{1, \dots, N_{Ph}\}} \left\{ \beta_j^{(k)} \right\} \delta \boldsymbol{\eta}^{(k)},$$

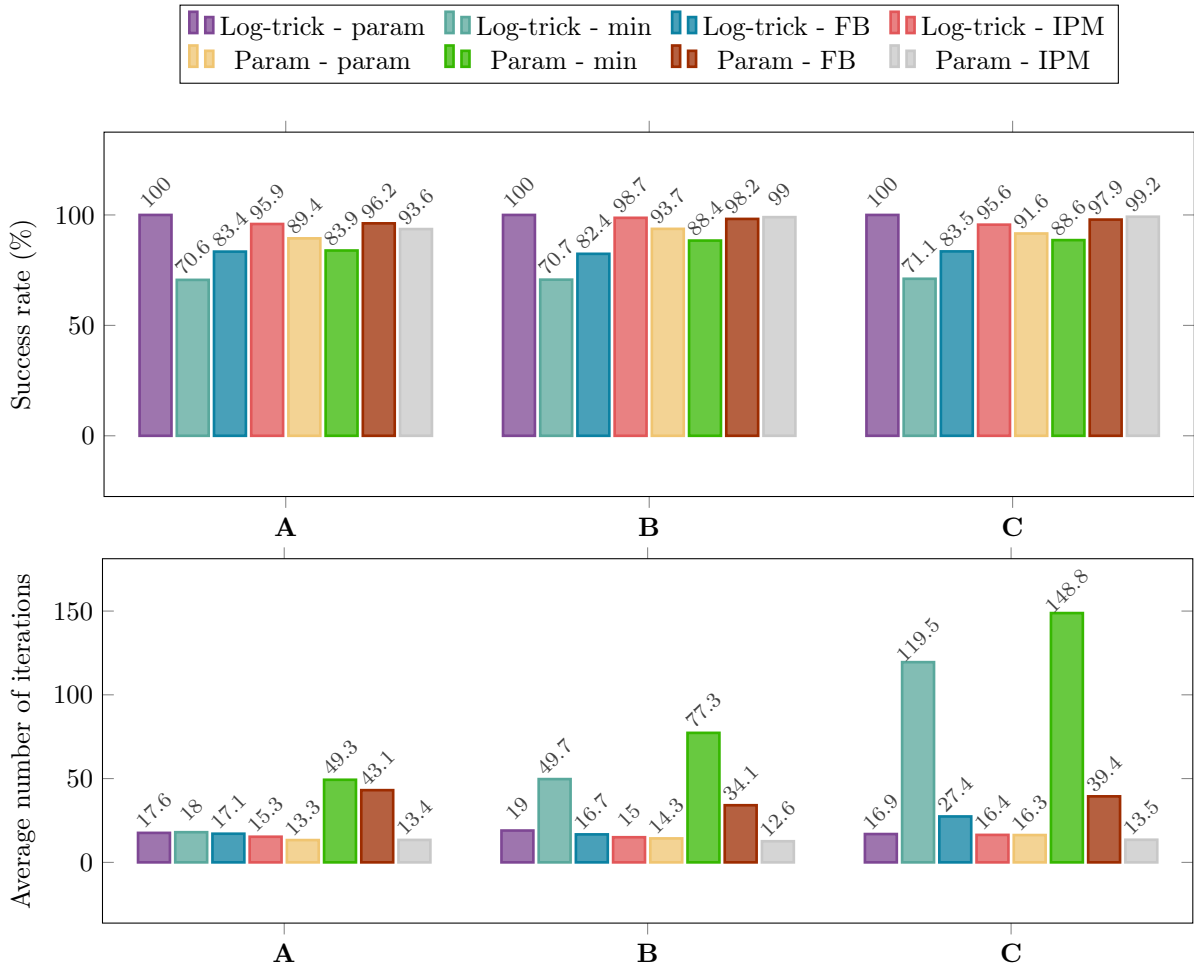
$$\beta_j^{(k)} = \begin{cases} -\eta_j^{(k)} / \delta \eta_j^{(k)} & \text{if } \text{sign} \left( \eta_j^{(k)} + \delta \eta_j^{(k)} \right) \neq \text{sign} \left( \eta_j^{(k)} \right), \\ 1 & \text{otherwise.} \end{cases}$$

The performance of this approach is illustrated in Figure 13. We observe a significant improvement in success rates of previously unstable methods. Regarding convergence speed, the log-trick combined with complementarity parametrization is now competitive with the fastest methods. However, both the min and FB approaches can exhibit slow convergence, particularly in configuration **C**, where all phases are present—the min method being especially affected.

Overall, for this small chemical system, the combination of the log-trick with complementarity parametrization appears to offer the most favorable basin of attraction.



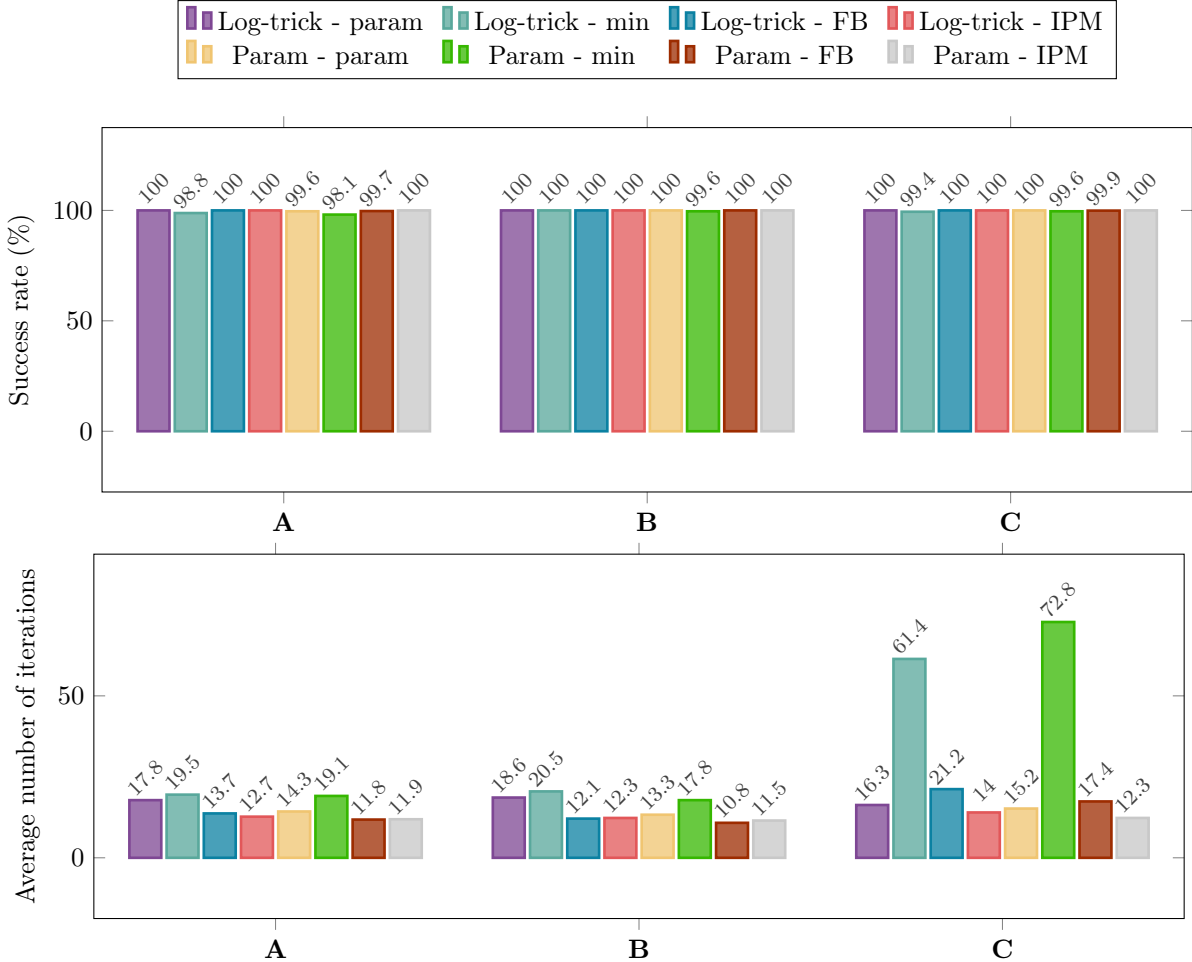
**Figure 12:** Performance metrics of CO<sub>2</sub> test case for a set of random initial guess.



**Figure 13:** Performance metrics of CO<sub>2</sub> test case for a set of random initial guess with a line search strategy.

**Sensitivity to right-hand side.** The results obtained for all methods using line search strategies are presented in Figure 14. All methods exhibit strong robustness with respect to variations in the right-hand side: the success rate is consistently close to or equal to 100%. In terms of convergence speed, the min function shows some difficulties, particularly in configuration C.

From these two sensitivity analyses, it is evident that incorporating a line search strategy is essential for ensuring robustness. Furthermore, the proposed complementarity-based parametrization proves to be competitive with, and in some cases superior to, classical methods.



**Figure 14:** Performance metrics of CO<sub>2</sub> test case for a set of random right-hand side with a line search strategy.

**Conclusion on CO<sub>2</sub> system and outlook.** This simplified CO<sub>2</sub> test case provides a controlled environment to validate the behavior of the different approaches. It highlights the significant role of the line search mechanism in enhancing robustness.

The combination of the log-trick with complementarity parametrization exhibits a wide basin of attraction and competitive convergence rates. In contrast, the FB and min functions may suffer from slower convergence in multi-phase contexts, especially without tailored damping strategies.

Overall, these preliminary results justify the use of advanced line search techniques and dedicated regularization mechanisms when tackling more challenging chemical systems. We now turn our attention to larger and more realistic systems, where phase behavior, nonlinearity, and degeneracy pose additional difficulties for robust convergence.

### 5.2.2 Multiphase Seawater: large-scale system for robustness evaluation

This second test case is composed of 72 species divided into 22 phases:

- 50 aqueous species in 1 phase;

Setting present phases	aqueous	mineral	gas
<b>A</b>	1	0	0
<b>B</b>	1	2	0
<b>C</b>	1	4	0
<b>D</b>	1	5	0
<b>E</b>	1	4	1

**Table 3:** Number of present phases at equilibrium for each setting.

- 20 minerals in 20 pure phases;
- 2 gaseous species in 1 phase.

It is a challenging test case that allows us to evaluate the limits of each method. There are no methods that can converge without a line search strategy. As for the previous test case, we have performed a study of the sensitivity to the initial guess and right-hand side. Tests were conducted on five chemical configurations (A to E), representing systems of increasing complexity:

**Sensitivity to initialization.** Figure 15 presents the performance for solving the *Seawater* case. Each configuration was tested over 1000 random initial guesses.

The results show that the parameterized formulation of the log significantly outperforms the *log-trick* in terms of success rate, especially when multiple mineral phases are present. While the *log-trick* performs well on simpler settings (A, B) with the complementarity parametrization method (>90%), it drastically fails in more complex cases (C to E), with success rates dropping below 40%. The FB and IPM strategies are almost entirely ineffective in this formulation. In contrast, the parameterized formulation of the log maintains high and stable success rates across all configurations when combined with either the complementarity parametrization (“param”) method (around 88%) or the FB method (approximately 80% across most settings, with a slight drop in E). The min method, however, is practically unusable: it fails to converge in the vast majority of simulations, resulting in a 0% success rate.

In terms of efficiency, measured by the average number of iterations, both “param” and FB methods in the parameterized formulation of the log are competitive. The “param” method is the fastest (30 to 50 iterations), while the FB method remains efficient with iteration counts between 40 and 80.

In summary, the combination of the parameterized formulation with the complementarity parametrization exhibits the highest overall performance.

**Convergence analysis with random perturbations on the right-hand side vector.** Figure 16 presents the performance when the right-hand side vector  $\mathbf{b}$  is randomly perturbed across 1000 simulations. Due to the nature of the perturbations, it is difficult to strictly control the number of phases present at equilibrium. However, the majority of simulations still fall within the five main chemical equilibrium configurations (A to E).

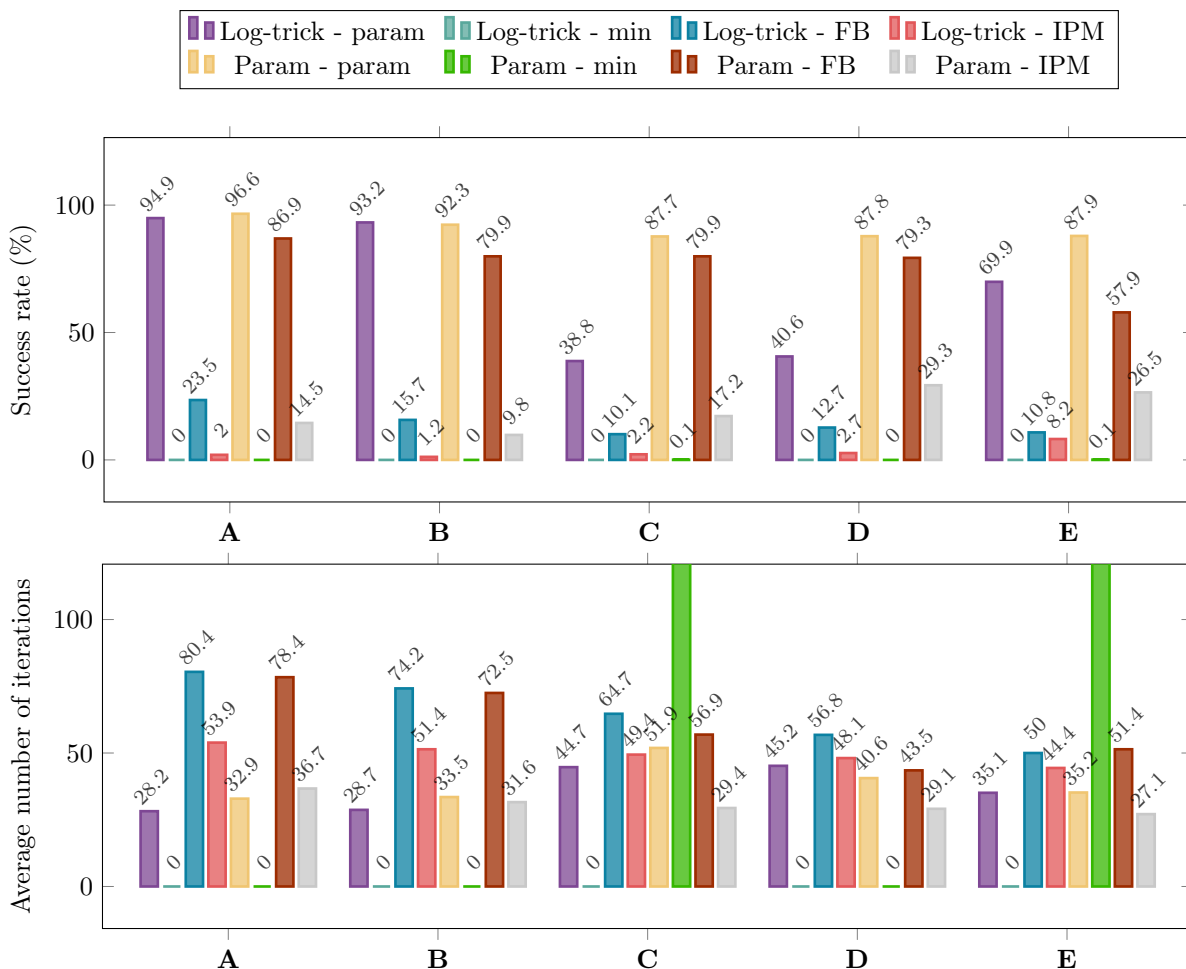
Overall, the success rates remain very high for the parameterized formulation of the log combined with the “param” method, consistently exceeding 90% even for the most complex configurations (C, D, and E). This robustness contrasts with the *log-trick* formulation, which exhibits a marked drop in success rates for complex systems, falling below 50% for configurations C and D.

The FB approach in the parameterized formulation also achieves very good convergence rates (ranging from 55.6% to 89.1%), showing notable improvement compared to the *log-trick* + FB results. The IPM consistently underperforms across all configurations, with very low success rates. The min method again fails to converge in all cases, confirming its unreliability under random perturbations of  $b$ .

Regarding iteration counts, the parameterized formulation combined with the “param” method shows stable and moderate iteration numbers (around 33 to 47), which are comparable to those observed with the *log-trick* + param method. The FB method requires a higher number of iterations on average but remains within reasonable bounds (approximately 42 to 110 iterations).

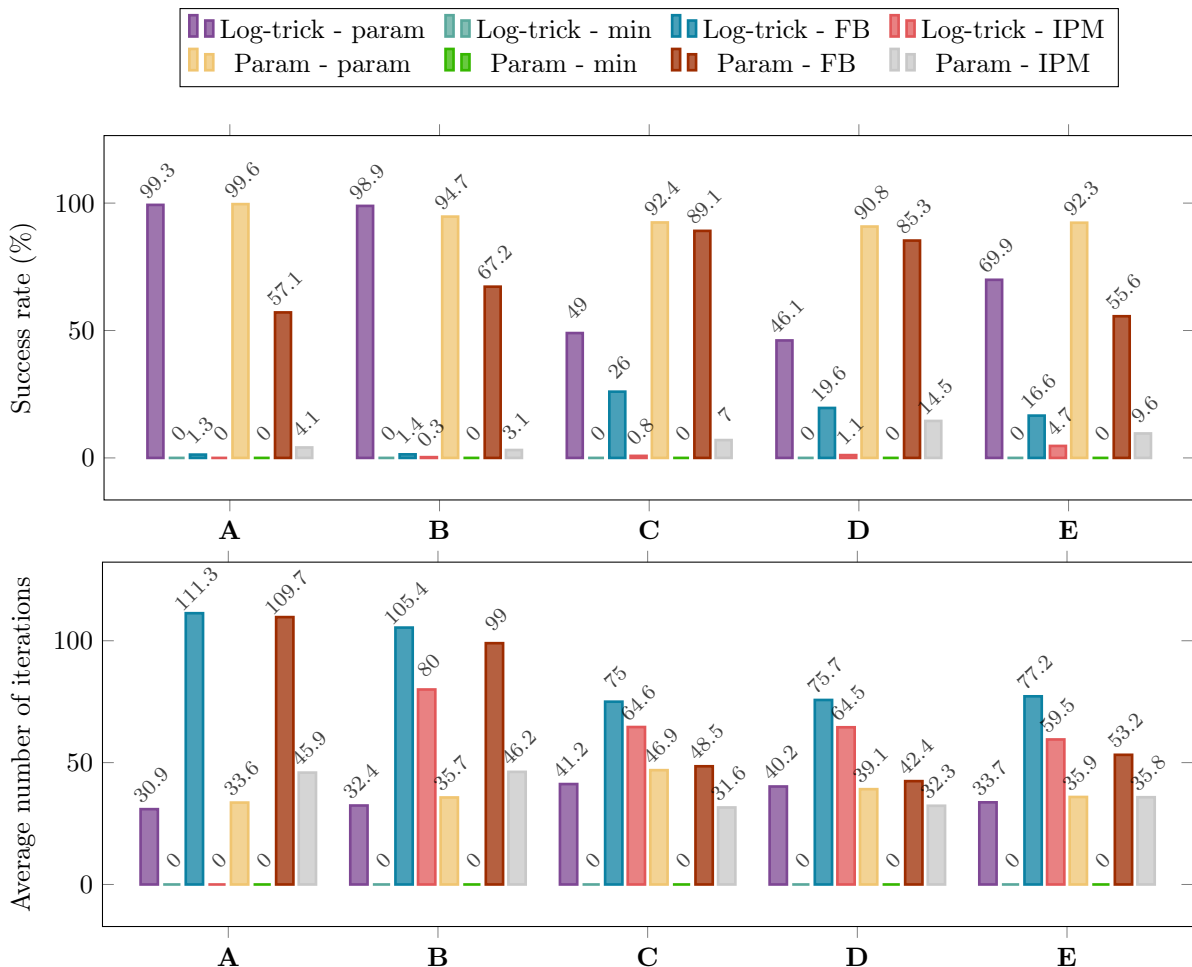
In conclusion, the parameterized formulation with the “param” approach once again demonstrates superior robustness and efficiency when facing perturbations in the right-hand side vector, confirming its suitability for solving challenging chemical equilibrium problems involving complementarity constraints.

**Conclusion on multiphase Seawater system.** Among all the tested strategies, the combination of the parameterized formulation of the log with the parameterized complementarity approach exhibits the



**Figure 15:** Performance metrics of *Seawater* test case for a set of random initial guess with a line search strategy.

highest overall performance. It achieves the best convergence rates across all chemical configurations (above 87% in every case), while also requiring fewer Newton iterations on average compared to the other methods. This clearly indicates that the proposed parametrized approach is both more robust and more efficient than traditional methods used in the literature. These results validate the relevance of the parameterized complementarity formulation as a promising alternative for solving equilibrium problems involving complementarity constraints.



**Figure 16:** Performance metrics of *Seawater* test case for a set of random right-hand side with a line search strategy.

## 6 Conclusion

In this paper, we addressed the complexities of multiphase equilibrium by introducing new conditions that guarantee uniqueness even when certain phases may vanish. Additionally, we established an equivalence between the minimization problem and a new set of algebraic equations based on extended mole fractions, leveraging concepts from the subdifferential of Gibbs energy.

This work presented a comprehensive numerical evaluation of Newton-based methods for solving multiphase chemical equilibrium problems under complementarity constraints. We proposed two main innovations: (i) a parameterized reformulation of the link between quantities and chemical potentials, and (ii) a novel parametrization of the complementarity condition.

Through two distinct robustness tests — one based on random initializations, and the other on random perturbations of the right-hand side vector — in a challenging chemical system composed of 72 species and 22 phases, we observed a clear performance advantage for the proposed parameterized approaches. This approach consistently achieved the highest success rates (often exceeding 90%) across all equilibrium configurations, while requiring moderate iteration counts. This approach outperformed all classical techniques, including the commonly used Newton-min, Fischer-Burmeister and interior point methods, which demonstrated either poor robustness or high computational cost.

The classical *log-trick* formulation exhibited strong sensitivity to initial conditions and problem configurations, and appears to be significantly less robust in complex multiphase systems. Similarly, the interior point method failed to deliver reliable convergence across all tests, regardless of the formulation.

These results validate the proposed parametrization strategy as a powerful and reliable alternative to existing methods for solving constrained nonlinear systems arising in chemical thermodynamics. By improving both the convergence rate and robustness, especially in challenging configurations, our methods represent a promising direction for future developments in numerical geochemistry.

## CRedit authorship contribution statement

**Maxime Jonval:** Methodology, Software, Validation, Formal analysis, Investigation, Writing - Original Draft, Writing - Review & Editing, Visualization. **Ibtihel Ben Gharbia:** Conceptualization, Methodology, Validation, Writing - Review & Editing, Supervision. **Clément Cancès:** Conceptualization, Methodology, Formal analysis, Writing - Review & Editing, Supervision. **Thibault Faney:** Conceptualization, Methodology, Validation, Writing - Review & Editing, Supervision. **Quang-Huy Tran:** Conceptualization, Methodology, Writing - Review & Editing, Supervision.

## Acknowledgement

This work was jointly supported by IFPEN and Inria. Maxime Jonval and Clément Cancès acknowledge support from the Labex CEMPI (ANR-11-LABX-0007-01), whereas Ibtihel Ben Gharbia, Clément Cancès and Quang Huy Tran acknowledge support from the project MATHSOUT of the PEPR Mathematics in Interaction (ANR-23-EXMA-0010) funded by the French National Research Agency.

## A Data for chemical systems

### A.1 The SiO<sub>2</sub> test case

The SiO<sub>2</sub> test case is composed of

$$\mathcal{C} = (\text{H}_2\text{O}, \text{H}^+, \text{SiO}_2(\text{aq}), \text{OH}^-), \mathcal{E} = (\text{H}, \text{O}, \text{Si}), \mathcal{R} = (\text{OH}^- = \text{H}_2\text{O} - \text{H}^+, \text{SiO}_2(\text{Quartz}) = \text{SiO}_2(\text{aq})).$$

### A.2 The CO<sub>2</sub> test case

The CO<sub>2</sub> system is composed of

$$\begin{aligned} \mathcal{C} &= (\text{H}_2\text{O}, \text{Ca}^{2+}, \text{CO}_2(\text{aq}), \text{H}^+, \text{OH}^-, \text{HCO}_3^-, \text{CO}_3^{2-}, \text{CaCO}_3(\text{calcite}), \text{CO}_2(\text{g}), \text{H}_2\text{O}(\text{g})), \\ \mathcal{E} &= (\text{O}, \text{Ca}, \text{C}, \text{H}), \\ \mathcal{R} &= (\text{OH}^- = \text{H}_2\text{O} - \text{H}^+; \quad \text{HCO}_3^- = \text{H}_2\text{O} + \text{CO}_2(\text{aq}) - \text{H}^+ \\ &\quad \text{CO}_3^{2-} = \text{H}_2\text{O} + \text{CO}_2(\text{aq}) - 2\text{H}^+; \quad \text{CaCO}_3(\text{calcite}) = \text{H}_2\text{O} + \text{Ca}^{2+} + \text{CO}_2(\text{aq}) - 2\text{H}^+ \\ &\quad \text{CO}_2(\text{g}) = \text{CO}_2(\text{aq}); \quad \text{H}_2\text{O}(\text{g}) = \text{H}_2\text{O}) \end{aligned}$$

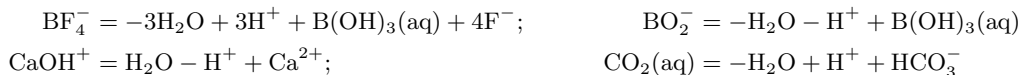
### A.3 The Seawater test case

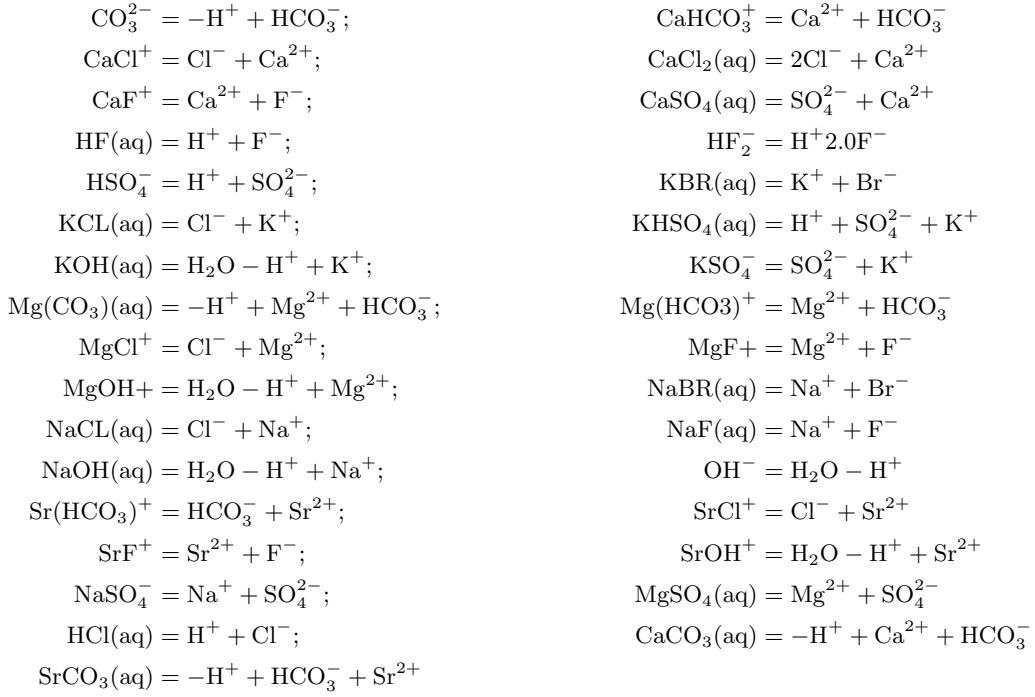
The *Seawater* test case is composed of

$$\begin{aligned} \mathcal{C} &= (\text{H}_2\text{O}(\text{l}), \text{H}^+, \text{Cl}^-, \text{Na}^+, \text{Mg}^{2+}, \text{SO}_4^{2-}, \text{Ca}^{2+}, \text{K}^+, \text{HCO}_3^-, \text{Br}^-, \\ &\quad \text{B}(\text{OH})_3(\text{aq}), \text{Sr}^{2+}, \text{F}^-, \text{BF}_4^-, \text{BO}_2^-, \text{CaOH}^+, \text{CO}_2(\text{aq}), \text{CO}_3^{2-}, \text{CaHCO}_3^+, \\ &\quad \text{CaCl}^+, \text{CaCl}_2(\text{aq}), \text{CaF}^+, \text{CaSO}_4(\text{aq}), \text{HF}(\text{aq}), \text{HF}_2^-, \text{HSO}_4^-, \text{KBR}(\text{aq}), \text{KCL}(\text{aq}), \\ &\quad \text{KHSO}_4(\text{aq}), \text{KOH}(\text{aq}), \text{KSO}_4^-, \text{Mg}(\text{CO}_3)(\text{aq}), \text{Mg}(\text{HCO}_3)^+, \text{MgCl}^+, \text{MgF}^+, \text{MgOH}^+, \text{NaBR}(\text{aq}), \\ &\quad \text{NaCL}(\text{aq}), \text{NaF}(\text{aq}), \text{NaOH}(\text{aq}), \text{OH}^-, \text{Sr}(\text{HCO}_3)^+, \text{SrCl}^+, \text{SrF}^+, \text{SrOH}^+, \text{NaSO}_4^-, \\ &\quad \text{MgSO}_4(\text{aq}), \text{HCl}(\text{aq}), \text{CaCO}_3(\text{aq}), \text{SrCO}_3(\text{aq}), \text{CaSO}_4(\text{anhydrite}), \text{CaCO}_3(\text{aragonite}), \\ &\quad \text{Mg}_2(\text{OH})_2(\text{CO}_3) \cdot 3\text{H}_2\text{O}(\text{artinite}), \text{Mg}(\text{OH})_2(\text{brucite}), \text{CaCO}_3(\text{calcite}), \text{SrSO}_4(\text{celestite}), \\ &\quad \text{CaMg}(\text{CO}_3)_2(\text{dolomite} - \text{dis}), \text{CaMg}(\text{CO}_3)_2(\text{dolomite} - \text{ord}), \text{CaF}_2(\text{fluorite}), \text{NaCl}(\text{halite}), \\ &\quad \text{CaMg}_3(\text{CO}_3)_4(\text{huntite}), \text{Mg}_5(\text{OH})_2(\text{CO}_3)_4 \cdot 4\text{H}_2\text{O}(\text{hydromagnesite}), \text{CaO}(\text{lime}), \\ &\quad \text{MgCO}_3(\text{magnesite}), \text{MgO}(\text{periclase}), \text{K}_2\text{O}(\text{potassium} - \text{oxide}), \text{Na}_2\text{O}(\text{sodium} - \text{oxide}), \\ &\quad \text{SrCO}_3(\text{strontianite}), \text{KCl}(\text{sylvite}), \text{MgCO}_3 \cdot 3\text{H}_2\text{O}(\text{nesquehonite}), \text{CO}_2(\text{g}), \text{H}_2\text{O}(\text{g})) \\ \mathcal{E} &= (\text{O}, \text{H}, \text{Cl}, \text{Na}, \text{Mg}, \text{S}, \text{Ca}, \text{K}, \text{C}, \text{Br}, \text{B}, \text{Sr}, \text{F}), \end{aligned}$$

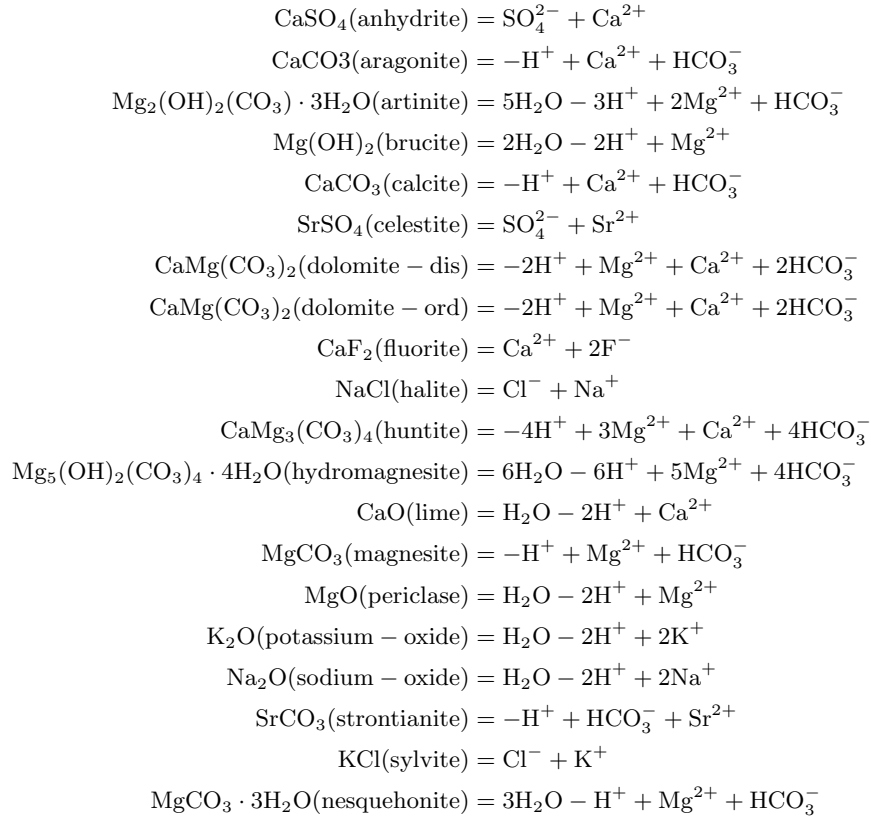
and the set  $\mathcal{R}$  composed of the reactions:

- aqueous species:

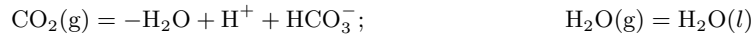




- pure mineral species:



- gaseous species:



## B Standard chemical potentials

The standard chemical potential  $\mu_i^\circ(P, T)$  of a species  $C_i$  for a constant pressure  $P$  and temperature  $T$  is calculated from the SUPCRT92 database [17].

**Table 4:** Standard chemical potentials at  $P = 1$  Bar and  $T = 298.15$  K.

Formula	$\mu_i^\circ(P, T)$	Formula	$\mu_i^\circ(P, T)$
H <sub>2</sub> O	-237138.9758928523	NaBR(aq)	-358192.19663367077
H <sup>+</sup>	0.0	NaCL(aq)	-388735.3787766348
Cl <sup>-</sup>	-131289.73255269116	NaF(aq)	-537936.7588824253
Na <sup>+</sup>	-261880.68093357567	NaOH(aq)	-417981.5046476963
Mg <sup>2+</sup>	-453984.7875582697	OH <sup>-</sup>	-157297.4420362631
SO <sub>4</sub> <sup>2-</sup>	-744458.9698933411	Sr(HCO <sub>3</sub> ) <sup>+</sup>	-1.157796221963139e6
Ca <sup>2+</sup>	-552789.9276571299	SrCl <sup>+</sup>	-693707.0438437797
K <sup>+</sup>	-282461.785083038	SrF <sup>+</sup>	-846381.1427032378
HCO <sub>3</sub> <sup>-</sup>	-586939.7841306848	SrOH <sup>+</sup>	-725087.037231101
Br <sup>-</sup>	-104056.09221446048	NaSO <sub>4</sub> <sup>-</sup>	-1.0103353876813356e6
B(OH) <sub>3</sub> (aq)	-968763.1798276871	MgSO <sub>4</sub> (aq)	-1.2111714801953817e6
Sr <sup>2+</sup>	-563835.6855778464	HCl(aq)	-127235.40484113985
F <sup>-</sup>	-281750.4962683066	CaCO <sub>3</sub> (aq)	-1.099764116286581e6
BF <sub>4</sub> <sup>-</sup>	-1.4869933045074963e6	SrCO <sub>3</sub> (aq)	-1.1081739548331397e6
BO <sub>2</sub> <sup>-</sup>	-678811.9589547777	CaSO <sub>4</sub> (Anhydrite)	-1.3218299449455e6
CaOH <sup>+</sup>	-716719.0378653706	CaCO <sub>3</sub> (Aragonite)	-1.12835344756575e6
CO <sub>2</sub> (aq)	-385973.95119834575	Mg <sub>2</sub> (OH) <sub>2</sub> (CO <sub>3</sub> ) <sub>2</sub> · 3H <sub>2</sub> O(Artinite)	-2.5686198744775e6
CO <sub>3</sub> <sup>2-</sup>	-527983.0213241152	Mg(OH) <sub>2</sub> (Brucite)	-835318.69698875
CaHCO <sub>3</sub> <sup>+</sup>	-1.145704464880173e6	CaCO <sub>3</sub> (Calcite)	-1.1291776990575502e6
CaCl <sup>+</sup>	-682410.2459630222	SrSO <sub>4</sub> (Celestite)	-1.34069978084625e6
CaCl <sub>2</sub> (aq)	-811695.8491350175	CaMg(CO <sub>3</sub> ) <sub>2</sub> (Dolomite – dis)	-2.1574917216637502e6
CaF <sup>+</sup>	-838431.5438473737	CaMg(CO <sub>3</sub> ) <sub>2</sub> (Dolomite – ord)	-2.1663074074905002e6
CaSO <sub>4</sub> (aq)	-1.3092988132949607e6	CaF <sub>2</sub> (Fluorite)	-1.17358246402515e6
HF(aq)	-299833.75289883884	NaCl(Halite)	-384120.431987875
HF <sub>2</sub> <sup>-</sup>	-578061.3350002527	CaMg <sub>3</sub> (CO <sub>3</sub> ) <sub>4</sub> (Huntite)	-4.2037057981325e6
HSO <sub>4</sub> <sup>-</sup>	-755755.7895325241	Mg <sub>5</sub> (OH) <sub>2</sub> (CO <sub>3</sub> ) <sub>4</sub> · 4H <sub>2</sub> O(Hydromagnesite)	-5.864657282973e6
KBR(aq)	-376601.8109260084	CaO(Lime)	-604027.2218463
KCL(aq)	-399279.0752045901	MgCO <sub>3</sub> (Magnesite)	-1.02783286346585e6
KHSO <sub>4</sub> (aq)	-1.0183854278434662e6	MgO(Periclase)	-569383.7028176
KOH(aq)	-437227.9151528981	K <sub>2</sub> O(Potassium – oxide)	-322402.2804475
KSO <sub>4</sub> <sup>-</sup>	-1.0319415534918394e6	Na <sub>2</sub> O(Sodium – oxide)	-376070.415242
Mg(CO <sub>3</sub> ) <sub>2</sub> (aq)	-998971.5788823194	SrCO <sub>3</sub> (Strontianite)	-1.15256625621825e6
Mg(HCO <sub>3</sub> ) <sup>+</sup>	-1.0468365642174695e6	KCl(Sylvite)	-408923.19198430004
MgCl <sup>+</sup>	-584504.6645023846	MgCO <sub>3</sub> · 3H <sub>2</sub> O(Nesquehonite)	-1.7239541270617498e6
MgF <sup>+</sup>	-743454.7620717564	CO <sub>2</sub> (g)	-394391.27240993205
MgOH <sup>+</sup>	-624482.7757953086	H <sub>2</sub> O(g)	-228164.33933913204

## C State of the art methods for the resolution

### C.1 The log-trick

As for the single-phase case, it could be interesting to test the robustness of the *log-trick* approach as a reference for comparing the methods. The idea of this approach is to make the following change of variable:

$$y = \ln \xi.$$

This ensures that the mole fractions are always positive and prevents issues with the Jacobian when  $\xi$  is small. Applying this change of variable to (12a)–(12c) yields:

$$\sum_{\alpha=1}^{N_{Ph}} s_{\alpha} \mathbf{A}^{\alpha} \exp(\mathbf{y}^{\alpha}) - b = 0, \quad (54a)$$

$$\mathbf{S}^T [\mu_i^{\circ}/(RT) + y_i]_{i=1, \dots, N} = \mathbf{0}, \quad (54b)$$

$$\langle \exp(\mathbf{y}^{\alpha}), \mathbf{1} \rangle + r_{\alpha} - 1 = 0, \quad (\alpha = 1, \dots, N_{Ph}), \quad (54c)$$

where  $\mathbf{exp}(\mathbf{y}) = (\exp(y_i))_{i=1, \dots, N}$ . The associated block in the Jacobian is written as:

$$\begin{bmatrix} s_1 \mathbf{A}^1 \text{diag}\{\mathbf{exp}(\mathbf{y}^1)\} & \cdots & s_{N_{Ph}} \mathbf{A}^{N_{Ph}} \text{diag}\{\mathbf{exp}(\mathbf{y}^{N_{Ph}})\} & \mathbf{A} \mathbf{exp}(\mathbf{y}) & \mathbf{0} \\ & & \mathbf{S}^T & \mathbf{0} & \mathbf{0} \\ & \mathbf{exp}(\mathbf{y}^1)^T & & & \\ & & \ddots & & \\ & & & \mathbf{0} & \mathbf{I}_{N_{Ph}} \\ & & & \mathbf{exp}(\mathbf{y}^{N_{Ph}})^T & \end{bmatrix}.$$

## C.2 Classical approaches to complementarity problem

Classical approaches to tackle the complementarity problem are typically categorized into two main classes of methods.

The first class consists of semismooth methods, where the complementarity problem is approached using a complementarity function. This function is Lipschitz continuous but may not be differentiable at specific points, particularly when  $s_\alpha = r_\alpha$ . However, by using the concept of subdifferentials, we can define a subgradient that allows the application of Newton's algorithm effectively. This approach facilitates the handling of non-differentiability while still leveraging the advantages of Newton's method. In Section C.3, we will present the minimum and Fischer-Burmeister complementarity functions, which are commonly used in semismooth approaches to effectively handle the complementarity conditions.

The second class focuses on smoothing techniques for the complementarity equations. In this approach, a regularization parameter is introduced to create a smooth approximation of the complementarity problem. This enables the use of classical Newton's method on the modified problem. As the iterations progress, the regularization parameter is gradually reduced to zero, guiding the solution back to the original complementarity problem as the iterates converge towards the solution. In Section C.4, we will explore the interior point method (IPM), a well-established technique for addressing the smoothing of complementarity equations.

## C.3 Semismooth methods

The treatment of the complementarity problem using semismooth methods involves the use of a complementarity function  $\Psi : \mathbb{R}^2 \rightarrow \mathbb{R}$  that satisfies specific conditions. Different functions lead to distinct smoothness and regularity characteristics, which are crucial for the effectiveness of the numerical techniques applied to solve nonlinear complementarity problems. To be a complementarity function,  $\Psi$  must satisfy the following condition:

$$\Psi(s_\alpha, r_\alpha) = 0 \quad \Leftrightarrow \quad s_\alpha r_\alpha = 0 \text{ and } s_\alpha \geq 0, r_\alpha \geq 0.$$

There are several options available for the function  $\Psi$ , but our study will concentrate on the following two functions:

- **The min function** [23]:

$$\Psi_{\min}(s_\alpha, r_\alpha) = \min(s_\alpha, r_\alpha),$$

which is Lipschitz continuous but lacks differentiability at the points where  $s_\alpha = r_\alpha$ . This function is depicted in Figure 17.

- **The Fischer-Burmeister function** [14]:

$$\Psi_{\text{FB}}(s_\alpha, r_\alpha) = s_\alpha + r_\alpha - \sqrt{s_\alpha^2 + r_\alpha^2},$$

which is differentiable everywhere except at the point where  $s_\alpha = r_\alpha = 0$ . This function is illustrated in Figure 18.

The non-differentiability of the min and Fischer-Burmeister (FB) functions leads to the use of semismooth Newton's method to solve the complementarity problem. In this approach, the Jacobian is replaced by an element of the Clarke subdifferential of the function for which we seek the root.

Let  $\mathcal{F} : \mathcal{X} \in \mathcal{D}_{\mathcal{F}} \rightarrow \mathbb{R}$  be a locally Lipschitz continuous function defined in a set  $\mathcal{D}_{\mathcal{F}}$ . By Rademacher's theorem [10, Section 3.4.1], such a function is continuously differentiable almost everywhere. Let  $\Omega_{\mathcal{F}}$  be the set where  $\mathcal{F}$  is Fréchet differentiable, *i.e.*, where  $\nabla \mathcal{F}$  exists. The Bouligand subdifferential  $\partial_{\text{B}} \mathcal{F}(\mathcal{X})$  of  $\mathcal{F}$  at  $\mathcal{X} \in \Omega_{\mathcal{F}}$  is defined as:

$$\partial_{\text{B}} \mathcal{F}(\mathcal{X}) := \left\{ \mathcal{J} \mid \exists (\mathcal{X}^{(k)})_{k \in \mathbb{N}} \subset \Omega_{\mathcal{F}}, \mathcal{X}^{(k)} \rightarrow \mathcal{X}, \nabla \mathcal{F}(\mathcal{X}^{(k)}) \rightarrow \mathcal{J} \right\}.$$

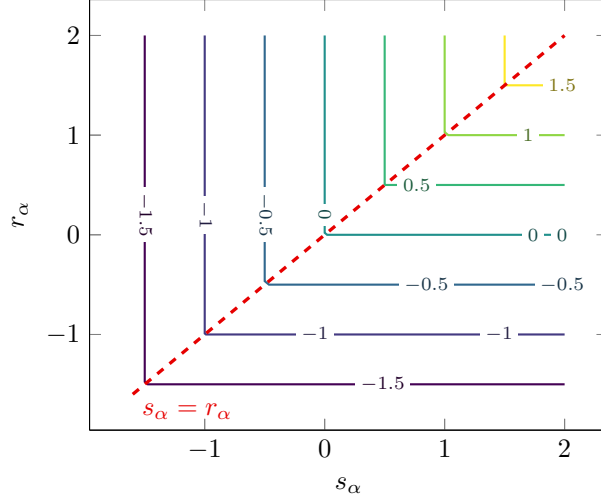


Figure 17: Contour plot of the function  $\Psi_{\min}$ .

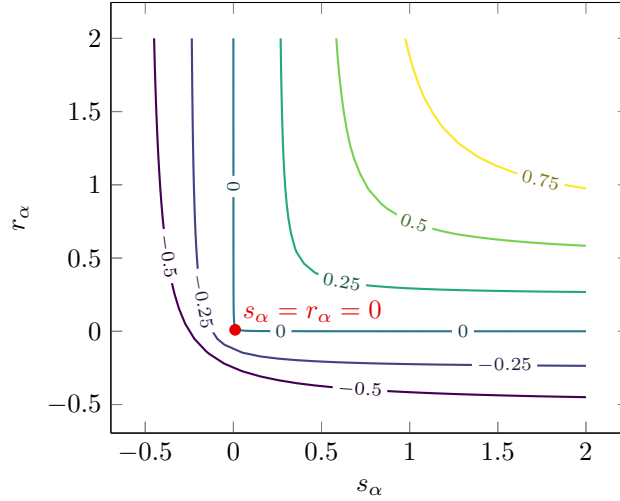


Figure 18: Contour plot of the function  $\Psi_{\text{FB}}$ .

The Clarke subdifferential [7] is then given by:

$$\partial_{\text{C}}\mathcal{F}(\boldsymbol{x}) = \text{conv}(\partial_{\text{B}}\mathcal{F}(\boldsymbol{x})),$$

where  $\text{conv}$  denotes the convex hull.

As a simple example, consider the absolute value function  $f(x) = |x|$  at  $x = 0$ . In this case,  $\partial_{\text{B}}f(0) = \{-1, 1\}$  and  $\partial_{\text{C}}f(0) = [-1, 1]$ .

The semismooth Newton's algorithm then start from an initial value  $\boldsymbol{x}^{(0)}$ , and builds a sequence  $(\boldsymbol{x}^{(k)})_{k>0}$  by solving the linear system [24]:

$$\mathcal{J}\delta\boldsymbol{x}^{(k)} = -\mathcal{F}(\boldsymbol{x}^{(k)}), \quad \mathcal{J} \in \partial_{\text{C}}\mathcal{F}(\boldsymbol{x}^{(k)}),$$

and updating the sequence as:

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + \delta\boldsymbol{x}^{(k)}.$$

The selection of  $\mathcal{J} \in \partial_{\text{C}}\mathcal{F}(\boldsymbol{x}^{(k)})$  when the latter is multivalued can have an impact on the performance of the method. Some authors [16] suggest choosing an element  $\mathcal{J}$  from the Bouligand subdifferential  $\partial_{\text{B}}\mathcal{F}(\boldsymbol{x}^{(k)})$  instead. In our case, the corresponding Bouligand subdifferentials for the min and FB functions are:

$$\begin{aligned} \partial_{\text{B}}\Psi_{\min}(s_{\alpha}, r_{\alpha})|_{s_{\alpha}=r_{\alpha}} &= \{(1, 0)^T, (0, 1)^T\} \\ \partial_{\text{B}}\Psi_{\text{FB}}(0, 0) &= \{(1, 1)^T - \mathbf{u} \mid \|\mathbf{u}\|_2 = 1\}. \end{aligned}$$

In our algorithms, the chosen subdifferentials are:

$$\begin{aligned}\partial_{\mathbf{B}}\Psi_{\min}(s_{\alpha}, r_{\alpha})|_{s_{\alpha}=r_{\alpha}} &= (1, 0)^T \\ \partial_{\mathbf{B}}\Psi_{\text{FB}}(0, 0) &= (1 - \sqrt{2}/2, 1 - \sqrt{2}/2)^T.\end{aligned}$$

However, these specific cases almost never occur in our numerical experiments.

## C.4 Smoothing methods

The central concept of smoothing techniques involves the introduction of a new parameter, denoted by  $\nu \in \mathbb{R}$ , known as the regularization parameter, into the complementarity equation. This parameter enables the "smoothing" of the equation, rendering it differentiable and thus suitable for classical optimization methods, such as the Newton method. In this context, we will focus on a specific class of smoothing methods that are part of interior point techniques [15]. Here, the complementarity equation (12d) is regularized as follows:

$$s_{\alpha}r_{\alpha} = \nu.$$

The system is then transformed into a new, smooth equation given by:

$$\mathcal{F}(\mathcal{X}; \nu) = \mathbf{0}. \quad (55)$$

The smooth equation can be effectively solved using the Newton method. During the Newton iterations, the regularization parameter is gradually decreased, allowing the smooth equation to converge toward the original complementarity equation. As the iterations progress and  $\nu$  approaches zero, the solution of the smooth equation converges to that of the original complementarity equation. The Newton step is defined by solving:

$$\nabla \mathcal{F}(\mathcal{X}^{(k)}; \nu^{(k)}) \delta \mathcal{X}^{(k)} = -\mathcal{F}(\mathcal{X}^{(k)}; \nu^{(k)}).$$

Subsequently,  $\nu$  is updated according to a predefined sequence:

$$\nu^{(k+1)} = \Theta(\nu^{(k)}),$$

which converges toward zero. In our study, we consider the following strategy for reducing  $\nu^{(k)}$  [32]:

$$\nu^{(k+1)} = \max \left( 10^{-32}, \min \left( \kappa_{\nu} \nu^{(k)}, \left( \nu^{(k)} \right)^{\theta_{\nu}} \right) \right),$$

with  $\kappa_{\nu} \in (0, 1)$  and  $\theta_{\nu} \in (1, 2)$ . At this stage, the complementarity equation is not fully satisfied, as we must ensure  $s_{\alpha}, r_{\alpha} \geq 0$ . While several strategies exist to maintain the positivity of these iterates, we present a method that treats  $s_{\alpha}$  and  $r_{\alpha}$  individually. The updates for  $\mathbf{s} = (s_{\alpha})_{\alpha=1, \dots, N_{Ph}}$  and  $\mathbf{r} = (r_{\alpha})_{\alpha=1, \dots, N_{Ph}}$  are as follows:

$$\begin{aligned}\mathbf{s}^{(k+1)} &= \mathbf{s}^{(k)} + \beta_{\mathbf{s}}^{(k)} \delta \mathbf{s}^{(k)}, \\ \mathbf{r}^{(k+1)} &= \mathbf{r}^{(k)} + \beta_{\mathbf{r}}^{(k)} \delta \mathbf{r}^{(k)}.\end{aligned}$$

Following the approach outlined in [32],  $\beta_{\mathbf{u}}$  for  $\mathbf{u} \in \{\mathbf{s}, \mathbf{r}\}$  is defined as:

$$\beta_{\mathbf{v}} = \max\{\beta \in (0, 1] \mid \mathbf{u} + \beta \delta \mathbf{u} \geq (1 - \gamma^{(k)}) \mathbf{u}\},$$

where  $\gamma^{(k)}$  is referred to as the fraction-to-the-boundary parameter in [32]:

$$\gamma^{(k)} = \max(\gamma_{\min}, 1 - \nu^{(k)}),$$

with  $\gamma_{\min} \in (0, 1)$ . We will refer to this update as the IPOPT line search.

Let us also mention an alternative method where  $\nu$  is considered as an unknown. This approach is named nonparametric IPM (NPIPM) and has been developed by Vu *et al.* [29, 30].

## References

- [1] V. Acary and B. Brogliato. *Numerical Methods for Nonsmooth Dynamical Systems*. Lecture Notes in Applied and Computational Mechanics. Springer Berlin, 2008.

- [2] I. Ben Gharbia. *Résolution de problèmes de complémentarité. : Application à un écoulement diphasique dans un milieu poreux*. PhD thesis, Université Paris Dauphine - Paris IX, 2012.
- [3] I. Ben Gharbia and E. Flauraud. Study of compositional multiphase flow formulation using complementarity conditions. *Oil Gas Sci. Technol. – Rev. IFP Energies nouvelles*, 74(43), 2019.
- [4] I. Ben Gharbia, M. Haddou, Q.-H. Tran, and D. T. S. Vu. An analysis of the unified formulation for the equilibrium problem of compositional multiphase mixtures. *ESAIM: M2AN*, 55(6):38, 2021.
- [5] I. Ben Gharbia and J. Jaffré. Gas phase appearance and disappearance as a problem with complementarity constraints. *Mathematics and Computers in Simulation*, 99:28–36, 2014.
- [6] K. Brenner and C. Cancès. Improving Newton’s method performance by parametrization: The case of the Richards equation. *SIAM Journal on Numerical Analysis*, 55(4):1760–1785, 2017.
- [7] F. H. Clarke. Generalized gradients and applications. *Transactions of the American Mathematical Society*, 205:247–262, 1975.
- [8] J. Coatléven and A. Michel. A successive substitution approach with embedded phase stability for simultaneous chemical and phase equilibrium calculations. *Computers & Chemical Engineering*, 168:108041, 2022.
- [9] K. H. Coats. An equation of state compositional model. *Society of Petroleum Engineers Journal*, 20(05):363–376, 10 1980.
- [10] Ş. Cobzaş, R. Miculescu, and A. Nicolae. *Lipschitz Functions*, volume 2241 of *Lecture Notes in Mathematics*. Springer Cham, 2019.
- [11] J. P. Dussault, M. Frappier, and J. Ch. Gilbert. Polyhedral Newton-min algorithms for complementarity problems. *Mathematical Programming*, May 2025.
- [12] F. Facchinei and J.-S. Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*, volume Volume I of *Springer Series in Operations Research and Financial Engineering*. Springer New York, 2003.
- [13] F. Facchinei and J.-S. Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*, volume Volume II of *Springer Series in Operations Research and Financial Engineering*. Springer New York, 2003.
- [14] A. Fischer. A special newton-type optimization method. *Optimization*, 24(3-4):269–284, 1992.
- [15] J. Gondzio. Interior point methods 25 years later. *European Journal of Operational Research*, 218(3):587–601, 2012.
- [16] A. F. Izmailov and M. V. Solodov. *Newton-Type Methods for Optimization and Variational Problems*. Springer Series in Operations Research and Financial Engineering. Springer Cham, 2014.
- [17] J. W. Johnson, E. H. Oelkers, and H. C. Helgeson. SUPCRT92: A software package for calculating the standard molal thermodynamic properties of minerals, gases, aqueous species, and reactions from 1 to 5000 bar and 0 to 1000°C. *Computers & Geosciences*, 18(7):899–947, 1992.
- [18] M. Jonval, I. Ben Gharbia, C. Cancès, T. Faney, and Q.-H. Tran. Parametrization and cartesian representation techniques for robust resolution of chemical equilibria. *Journal of Computational Physics*, 522:113596, 2025.
- [19] A. Lauser, C. Hager, R. Helmig, and B. Wohlmuth. A new approach for phase transitions in miscible multi-phase flow in porous media. *Advances in Water Resources*, 34(8):957–966, 2011.
- [20] A. M. M. Leal, D. A. Kulik, W. R. Smith, and M. O. Saar. An overview of computational methods for chemical equilibrium and kinetic calculations for geochemical and reactive transport modeling. *Pure and Applied Chemistry*, 89(5):597–643, 2017.
- [21] M. L. Michelsen. The isothermal flash problem. part ii. phase-split calculation. *Fluid Phase Equilibria*, 9(1):21–40, 1982.

- [22] D. K. Nordstrom and K. M. Campbell. Modeling low-temperature geochemical processes. In *Treatise on Geochemistry*, pages 27–68. Elsevier, 2014.
- [23] J.-S. Pang. Newton’s method for b-differentiable equations. *Mathematics of Operations Research*, 15(2):311–341, 1990.
- [24] L. Qi and J. Sun. A nonsmooth version of newton’s method. *Mathematical Programming*, 58:353–367, 1993.
- [25] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [26] N. Z. Shapiro and L. S. Shapley. Mass action laws and the Gibbs free energy function. *Journal of the Society for Industrial and Applied Mathematics*, 13(2):353–375, 1965.
- [27] J. V. Smith, R. W. Missen, and W. R. Smith. General optimality criteria for multiphase multireaction chemical equilibrium. *AIChE Journal*, 39(4):707–710, 1993.
- [28] W. R. Smith and R. W. Missen. *Chemical reaction equilibrium analysis: Theory and algorithms*. John Wiley & Sons, New York, 1982.
- [29] D. T. S. Vu. Numerical resolution of algebraic systems with complementarity conditions: Application to the thermodynamics of compositional multiphase mixtures. *Université Paris-Saclay*, 2020.
- [30] D. T. S. Vu, I. Ben Gharbia, M. Haddou, and Q.-H. Tran. A new approach for solving nonlinear algebraic systems with complementarity conditions. application to compositional multiphase equilibrium problems. *Mathematics and Computers in Simulation*, 190:1243–1274, 2021.
- [31] T.J. Wolery. EQ3NR, a computer program for geochemical aqueous speciation-solubility calculations: Theoretical manual, user’s guide, and related documentation (Version 7.0); Part 3. Technical Report UCRL-MA-110662-Pt.3, 138643, Lawrence Livermore National Laboratory, 1992.
- [32] A. Wächter and L. T. Biegler. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106(1):25–57, 2006.