

FRIEDRICH-ALEXANDER-UNIVERSITÄT ERLANGEN-NÜRNBERG
INSTITUT FÜR ANGEWANDTE MATHEMATIK

A general reduction scheme for reactive
transport in porous media

by

J. Hoffmann, S. Krättele & P. Knabner

No. 353

2012



Institut für Angewandte Mathematik
Martensstraße 3
D-91058 Erlangen

Lehrstuhl AM I:

Tel.: ++49/9131/85-27015

Fax: ++49/9131/85-27670

Lehrstuhl AM II:

Tel.: ++49/9131/85-27510

Fax: ++49/9131/85-28126

Lehrstuhl AM III:

Tel.: ++49/9131/85-25226

Fax: ++49/9131/85-25228

internet: <http://www.am.uni-erlangen.de>

ISSN 1435-5833, 2012, No. 353

A general reduction scheme for reactive transport in porous media

Joachim Hoffmann, Serge Kräutle, Peter Knabner¹

Abstract We present a method to transform the governing equations of multi-species reactive transport in porous media. The reformulation leads to a smaller problem size by decoupling of equations and by elimination of unknowns, which increases the efficiency of numerical simulations. The reformulation presented here is a generalization of earlier works. In fact, a whole class of transformations is now presented. This class is parametrized by the choice of certain transformation matrices. For specific choices, some known formulations of reactive transport can be retrieved. Hence, the software based on the presented transformation can be used to obtain efficiency comparisons of different solution approaches. For our efficiency tests we use the MoMaS Benchmark Problem on Reactive Transport.

Keywords reactive transport, porous media, reduction of problem size, numerical simulation, global implicit approach

MSC(2010) 35K57, 35K58, 65Y20

1 Introduction

This work is devoted to efficient numerical simulations of large-scale reactive transport problems in porous media. We present a reformulation of the given system of partial differential equations (PDEs), ordinary differential equations (ODEs), algebraic equilibrium equations (AEs) and inequalities that govern multicomponent reactive transport problems. The reformulation uses linear transformations of the equations and variables. The size of the governing system of equations is reduced by decoupling of some linear PDEs and by implicit elimination of local equations (i.e., AEs, ODEs) and unknowns. The work presented here is based on earlier versions of the reduction scheme which are proposed in [13, 14, 11].

Let us emphasize that our method is a global implicit approach (GIA) and does neither use operator splitting nor simplifications of the model. Operator splitting schemes usually decompose a time step into a transport problem and a reaction problem. There are two main types, the sequential non-iterative approach (SNIA) and the sequential iterative approach (SIA).

¹Email: kraeutle/knabner@am.uni-erlangen.de

University of Erlangen-Nuremberg, Department of Mathematics, Cauerstr. 11, 91058 Erlangen, Germany

Submitted: Feb 21, 2012.

For example the software SPECY uses a non-iterative operator splitting scheme (see [3]). The drawback of SNIA is the introduction of splitting errors (see e.g. [6]). To circumvent this problem, iterative operator splitting schemes can be used. For example the software HYTEC uses an iterative operator splitting scheme (see [22]). However, iterative operator splitting schemes may need many iteration steps and small time steps to converge in chemically difficult cases which can impact the effectency (see e.g. [20]). An advantage of splitting schemes seems that the implementation might be easier compared to a reformulation scheme. One GIA method, which is known for several years, is the direct substitutional approach (DSA). For example the code MIN3P (see [17]) uses this method. One disadvantage of DSA is that it leads to a nonlinear system which might be difficult to solve numerically. A recently proposed global implicit approach is to use a differential algebraic equation (DAE) solver (see [7, 8]). Here the transport equations, the mass balance equations and the equations describing the chemical equilibrium are solved in one very large system of equations. The disadvantage is that this approach seems to lead to larger computation times (see [4]).

Besides the presentation of the reformulation technique, the other aim of this paper is to give a efficiency comparison to other methods, based on a demanding benchmark problem [5], which is characterized by strongly nonlinear reactions, variation of concentrations and scales of many magnitudes, and two space dimensions.

Compared to earlier presentations [14, 9, 10], the reformulation presented here is more general. In fact, a whole class of transformations is described. In this class of transformations, not only [14, 9, 10], but also the Morel formulation [5, 18], which is a common standard for such problems, is found as a representative.

Hence, with one implementation comparisons between the reduction scheme and the Morel formulation are possible. That excludes differences in the CPU time due to the specific implementation and specific computer systems, etc. We hope that this contributes to unbiased comparisons of the different methods.

The paper is organized as follows. In Sec. 2 we present the model. It includes mobile species and immobile species such as sorbed species and minerals, and it includes both kinetic and equilibrium reactions. In Sec. 3 we develop the reformulation and size reduction technique. Sec. 4 is devoted to the question if the substitution of local equations into the remaining equations, which is proposed in Sec. 3, is always possible. Sec. 5 shows that both the formulation [14] and the Morel formulation [5] are special cases of the presented general scheme. In Sec. 6 we give numerical results for the

schemes for the MoMaS benchmark [5]. These results can be seen as an extension of the results presented in [4, 10].

2 Mathematical model

Our model of multicomponent reactive transport covers both kinetic and equilibrium reactions, and both mobile and immobile species. The details are as follows.

2.1 Chemical reactions

Kinetic reactions according to law of mass action

Concerning kinetic reactions according to the law of mass action the rate term is given by the difference of the forward and the backward reaction rate (see e.g. [2])

$$r_{kin,j}(\mathbf{c}, \bar{\mathbf{c}}) = k_{f,j} \prod_{\substack{i=1 \\ s_{ij}<0}}^{I+\bar{I}} c_i^{-s_{ij}} - k_{b,j} \prod_{\substack{i=1 \\ s_{ij}>0}}^{I+\bar{I}} c_i^{+s_{ij}} \quad (1)$$

of the chemical reaction. The number of reactions of this type is denoted by J_{kin} . The number of mobile and immobile species is denoted by I, \bar{I} , respectively. In the following we will frequently mark immobile species by a bar (' \bar{c}_i ').

Another kind of kinetic reactions are biodegradation reactions. They can be described with help of the Monod model. A presentation of the model and how to apply the reduction mechanism to that type of reactions can be found in [13, Sec. 5].

Equilibrium reactions according to law of mass action

If the j -th equilibrium reaction can be described by the law of mass action then the j -th equilibrium condition reads

$$\phi_j(\mathbf{c}, \bar{\mathbf{c}}) := -\ln(K_j) + \sum_{i=1}^{I+\bar{I}} s_{ij} \ln(c_i) = 0. \quad (2)$$

In case of equilibrium reactions the corresponding reaction rate $r_{eq,j}$ is not known. So $r_{eq,j}$ gets an additional unknown and the equilibrium condition (2) is added to the system of PDEs and ODEs as an additional equation.

Equilibrium minerals

If the j -th equilibrium reaction is a mineral reaction we define

$$\psi_j(\mathbf{c}) := -\ln(K_j) + \sum_{i=1}^I s_{ij} \ln(c_i). \quad (3)$$

For a mineral reaction the equilibrium condition consisting of equations and inequalities reads (e.g. [12])

$$\left(\psi_j(\mathbf{c}) = 0 \wedge \bar{c}_{min,j} \geq 0 \right) \vee \left(\psi_j(\mathbf{c}) > 0 \wedge \bar{c}_{min,j} = 0 \right) \quad (4)$$

where $\bar{c}_{min,j}$ denotes the concentration of that mineral taking part in the j -th equilibrium reaction. The left part of (4) describes saturation, the right part describes undersaturation of the fluid with respect to mineral j .

The equilibrium condition (4) can be expressed by

$$\phi_j(\mathbf{c}, \bar{\mathbf{c}}_{min}) := \min \{ \psi_j(\mathbf{c}), \bar{c}_{min,j} \} = 0. \quad (5)$$

The equations (2) and (5) are subsumed to $\boldsymbol{\phi}(\mathbf{c}, \bar{\mathbf{c}}) = \mathbf{0}$ and added to the system of PDEs and ODEs, see next section.

2.2 Reactive transport

In matrix notation the system of equations reads

$$\begin{aligned} \partial_t(\theta \mathbf{c}) + L\mathbf{c} &= \theta \mathbf{S}_{1,kin} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) + \theta \mathbf{S}_{1,eq} \mathbf{r}_{eq} \\ \partial_t(\theta \bar{\mathbf{c}}) &= \theta \mathbf{S}_{2,kin} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) + \theta \mathbf{S}_{2,eq} \mathbf{r}_{eq} \\ \boldsymbol{\phi}(\mathbf{c}, \bar{\mathbf{c}}) &= \mathbf{0} \end{aligned} \quad (6)$$

with the Darcy velocity \mathbf{q} , the linear transport operator $L_i u_i := -\nabla \cdot (\mathbf{D}_i \nabla u_i - \mathbf{q} u_i)$, $L\mathbf{u} = \text{diag}(L_i u_i)$, and the stoichiometric matrix

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_{1,eq} & \mathbf{S}_{1,kin} \\ \mathbf{S}_{2,eq} & \mathbf{S}_{2,kin} \end{pmatrix}.$$

Some assumptions on the shape of the blocks of this matrix are given in the next section.

The Scheidegger diffusion/dispersion tensor \mathbf{D}_i (see [21]) is used to describe diffusion and dispersion

$$\mathbf{D}_i = -(\theta d_{diff,i} + \beta_t |\mathbf{q}|) \mathbf{I} + (\beta_l - \beta_t) \frac{\mathbf{q} \mathbf{q}^T}{|\mathbf{q}|}. \quad (7)$$

This system of equations is widely used for modelling reactive transport (see, e.g., [19, eq. (19),(1)]).

3 The general reduction mechanism

The reduction scheme presented here is an extension to the reduction scheme described in [13], [14], [11] and [10]. To apply the reduction scheme the assumption that the diffusion coefficient $d_{diff,i}$ is the same for all species is needed. This assumption is justified because the molecular diffusion is usually small compared to the mechanic dispersion.

We arrange the equilibrium reactions in the following order. First we take the reactions in which only mobile species take part, then the equilibrium sorption reactions, i.e., heterogeneous without mineral species, and at last the equilibrium mineral reactions:

$$\begin{aligned} \mathbf{S}_{1,eq} &= (\mathbf{S}_{1,mob} \quad \mathbf{S}_{1,sorp} \quad \mathbf{S}_{1,min}), \\ \mathbf{S}_{2,eq} &= \begin{pmatrix} \mathbf{0} & \mathbf{S}_{2,sorp} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{J_{min}} \end{pmatrix} \end{aligned} \quad (8)$$

$\mathbf{I}_{J_{min}}$ denotes the identity matrix of the size J_{min} . The number of the reactions of each type is denoted by J_{mob} , J_{sorp} and J_{min} , respectively. Then we define the abbreviation $\mathbf{S}_{1,het} := (\mathbf{S}_{1,sorp} \quad \mathbf{S}_{1,min})$. Hence, we get

$$\mathbf{S}_{1,eq} = (\mathbf{S}_{1,mob} \quad \mathbf{S}_{1,het}) \quad (9)$$

Also the immobile species are arranged in a specific order. We take the nonminerals first and then the minerals $\bar{\mathbf{c}} = (\bar{\mathbf{c}}_{nmin}, \bar{\mathbf{c}}_{min})$. The number of nonminerals² is named \bar{J}_{nmin} .

We assume that both the columns of

$$(\mathbf{S}_{1,mob} \quad \mathbf{S}_{1,min}) \text{ and of } \mathbf{S}_{2,sorp} \quad (10)$$

are linear independent. Furthermore we assume³ that there is a subset of columns $\mathbf{S}_{1,sorp}^*$ of the columns of $\mathbf{S}_{1,sorp}$ and a subset of columns $\mathbf{S}_{1,kin}^*$ of the columns of $\mathbf{S}_{1,kin}$ in such a way that the columns of the matrix $(\mathbf{S}_{1,mob} \quad \mathbf{S}_{1,sorp}^* \quad \mathbf{S}_{1,min} \quad \mathbf{S}_{1,kin}^*)$ are a maximal system of linear independent columns of matrix \mathbf{S}_1 and that all the remaining columns of block $\mathbf{S}_{1,sorp}$ are linear combinations of the columns of $\mathbf{S}_{1,sorp}^*$, $\mathbf{S}_{1,min}$, $\mathbf{S}_{1,kin}^*$ (but not of $\mathbf{S}_{1,mob}$).

²In this work nonminerals always denote all *immobile* species that are not a mineral.

³This assumption is rather technical. It could be replaced by the simpler assumption that all columns of matrix $(\mathbf{S}_{1,mob} \quad \mathbf{S}_{1,sorp} \quad \mathbf{S}_{1,min})$ are linear independent, as it is done in [11, p.131]. Hence, the assumption after (10) is a relaxation of the assumption made in [11].

Now let us incorporate also the kinetic reactions. The matrices \mathbf{S}_1 and \mathbf{S}_2 , containing all entries of \mathbf{S} connected to mobile and immobile species, respectively, are of the form:

$$\mathbf{S}_1 = \begin{pmatrix} \mathbf{S}_{1,eq} & \mathbf{S}_{1,kin} \end{pmatrix} = \begin{pmatrix} \mathbf{S}_{1,mob} & \mathbf{S}_{1,het} & \mathbf{S}_{1,kin} \end{pmatrix} \quad (11)$$

$$\mathbf{S}_2 = \begin{pmatrix} \mathbf{S}_{2,eq} & \mathbf{S}_{2,kin} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{S}_{2,sorp} & \mathbf{0} & \tilde{\mathbf{S}}_{2,kin} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{J_{min}} & \mathbf{0} \end{pmatrix}$$

Note that this implies that we assume that an equilibrium mineral does not participate in any kinetic reaction. So the stoichiometric coefficients in $\mathbf{S}_{2,kin}$ connected to minerals are zero, and it is possible to express $\mathbf{S}_{2,kin}$ by $\tilde{\mathbf{S}}_{2,kin}$ and a zero block.

Recall that our system (6) reads

$$\partial_t(\theta\mathbf{c}) + L\mathbf{c} = \theta\mathbf{S}_1 \begin{pmatrix} \mathbf{r}_{eq} \\ \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix} \quad (12)$$

$$\partial_t(\theta\bar{\mathbf{c}}) = \theta\mathbf{S}_2 \begin{pmatrix} \mathbf{r}_{eq} \\ \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix} \quad (13)$$

$$\phi(\mathbf{c}, \bar{\mathbf{c}}) = \mathbf{0} \quad (14)$$

with the block structure (11) and equilibrium conditions (14) from (2), (5).

In order to prepare the reformulation, we choose matrices \mathbf{S}_1^* and \mathbf{S}_2^* that contain a maximal system of linear independent columns of \mathbf{S}_1 and \mathbf{S}_2 , respectively. Because of the linear independence assumption (see (10)) it is always possible to choose \mathbf{S}_1^* and \mathbf{S}_2^* such that the matrices have the form

$$\begin{aligned} \mathbf{S}_1^* &= \begin{pmatrix} \mathbf{S}_{1,mob} & \mathbf{S}_{1,het}^* & \mathbf{S}_{1,kin}^* \end{pmatrix} \in \mathbb{R}^{I \times (J_{mob} + J_{1,het}^* + J_{1,kin}^*)} \\ \mathbf{S}_2^* &= \begin{pmatrix} \mathbf{S}_{2,sorp} & \mathbf{0} & \mathbf{S}_{2,kin}^* \\ \mathbf{0} & \mathbf{I}_{J_{min}} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{\bar{I} \times (J_{sorp} + J_{min} + J_{2,kin}^*)} \end{aligned} \quad (15)$$

where $\mathbf{S}_{1,het}^*$, $\mathbf{S}_{1,kin}^*$ and $\mathbf{S}_{2,kin}^*$ consist of some of the columns of $\mathbf{S}_{1,het}$, $\mathbf{S}_{1,kin}$ and $\tilde{\mathbf{S}}_{2,kin}$, respectively. The number of columns of $\mathbf{S}_{1,het}^*$, $\mathbf{S}_{1,kin}^*$ and $\mathbf{S}_{2,kin}^*$ is denoted $J_{1,het}^*$, $J_{1,kin}^*$ and $J_{2,kin}^*$, respectively.

Since matrix \mathbf{S}_i^* consists of linear independent columns, there is always a matrix \mathbf{A}_i such that

$$\mathbf{S}_i = \mathbf{S}_i^* \mathbf{A}_i \quad i = 1, 2. \quad (16)$$

With the block structure from (11) for $\mathbf{S}_1, \mathbf{S}_2$ and the block structure from (15) for $\mathbf{S}_1^*, \mathbf{S}_2^*$, and with the assumptions made in (10) and thereafter, we get (similar as in [14]) for \mathbf{A}_1 and \mathbf{A}_2 the block structure

$$\begin{aligned} \mathbf{A}_1 &= \begin{pmatrix} \mathbf{I}_{J_{mob}} & \mathbf{0} & \mathbf{A}_{1,mob} \\ \mathbf{0} & \mathbf{A}_{1,het} & \tilde{\mathbf{A}}_{1,kin} \end{pmatrix} \\ &\in \mathbb{R}^{(J_{mob}+J_{1,het}^*+J_{1,kin}^*) \times (J_{mob}+J_{het}+J_{kin})} \\ \mathbf{A}_2 &= \begin{pmatrix} \mathbf{0} & \mathbf{I}_{J_{sorp}} & \mathbf{0} & \mathbf{A}_{2,sorp} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{J_{min}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{A}_{2,kin} \end{pmatrix} \\ &\in \mathbb{R}^{(J_{sorp}+J_{min}+J_{2,kin}^*) \times (J_{mob}+J_{sorp}+J_{min}+J_{kin})}, \end{aligned} \quad (17)$$

where $J_{mob} + J_{sorp} + J_{min} + J_{kin} = J_{mob} + J_{het} + J_{kin} = J$ is the total number of reactions.

Transformation of the PDE block

To derive the reduction scheme the block of the PDEs (12) is multiplied with some matrices $\tilde{\mathbf{S}}_1^{\perp T}$ and \mathbf{C}_1^T . In order to allow the user some freedom of choice, we do not prescribe specific matrices $\tilde{\mathbf{S}}_1^{\perp}$ and \mathbf{C}_1 . Instead, we give the following five conditions which have to be met by the matrices $\tilde{\mathbf{S}}_1^{\perp}, \mathbf{C}_1$ to be chosen. The conditions are formulated as general as possible. Note that two different specific choices of the matrices are presented in Sec. 5. The conditions are:

1. The number of rows of both matrices $\tilde{\mathbf{S}}_1^{\perp}, \mathbf{C}_1$ is I
2. $\text{span}\{\tilde{\mathbf{S}}_1^{\perp}, \mathbf{C}_1\} = \mathbb{R}^I$
3. All columns of $\tilde{\mathbf{S}}_1^{\perp}, \mathbf{C}_1$ are linear independent
4. All columns of $\tilde{\mathbf{S}}_1^{\perp}$ are orthogonal to all columns of \mathbf{S}_1^* (but it is not necessary that $\tilde{\mathbf{S}}_1^{\perp}$ is a maximal system of linear independent vectors that are orthogonal to \mathbf{S}_1^*): $\tilde{\mathbf{S}}_1^{\perp T} \mathbf{S}_1^* = \mathbf{0}$
5. There is a matrix \mathbf{D}_1 such that

$$\mathbf{C}_1^T \mathbf{S}_1^* = \begin{pmatrix} \mathbf{I}_{J_{mob}} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_1 \end{pmatrix} \quad (18)$$

(the block \mathbf{D}_1 may be rectangular)⁴

Since the matrix \mathbf{S}_1^* consists of $J_{mob} + J_{1,het}^* + J_{1,kin}^*$ columns, and 4. holds, the number of columns of $\tilde{\mathbf{S}}_1^\perp$ is less than or equal to $I - (J_{mob} + J_{1,het}^* + J_{1,kin}^*)$. Let N_* be defined by postulating that the number of columns of $\tilde{\mathbf{S}}_1^\perp$ is $I - J_{mob} - N_*$. It immediately follows that $J_{1,het}^* + J_{1,kin}^* \leq N_* \leq I - J_{mob}$. From 1., 2. and 3. it follows that the size of the matrix \mathbf{C}_1 is $I \times (J_{mob} + N_*)$. Hence, the product (18) has the size $(J_{mob} + N_*) \times (J_{mob} + J_{1,het}^* + J_{1,kin}^*)$. As a consequence, \mathbf{D}_1 has the size $N_* \times (J_{1,het}^* + J_{1,kin}^*)$.

Multiplying the block of the PDEs (12) by $\tilde{\mathbf{S}}_1^{\perp T}$ and by \mathbf{C}_1^T , and exploiting condition 4, we obtain the two equations

$$\begin{aligned} \partial_t \left(\theta \tilde{\mathbf{S}}_1^{\perp T} \mathbf{c} \right) + L \tilde{\mathbf{S}}_1^{\perp T} \mathbf{c} &= \mathbf{0} \\ \partial_t (\theta \mathbf{C}_1^T \mathbf{c}) + L \mathbf{C}_1^T \mathbf{c} &= \theta \mathbf{C}_1^T \mathbf{S}_1^* \mathbf{A}_1 \begin{pmatrix} \mathbf{r}_{eq} \\ \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix}. \end{aligned}$$

Using condition 5. and the structure of \mathbf{A}_1 (see (17)) it follows that

$$\begin{aligned} \partial_t \left(\theta \tilde{\mathbf{S}}_1^{\perp T} \mathbf{c} \right) + L \tilde{\mathbf{S}}_1^{\perp T} \mathbf{c} &= \mathbf{0}, \\ \partial_t (\theta \mathbf{C}_1^T \mathbf{c}) + L \mathbf{C}_1^T \mathbf{c} &= \theta \begin{pmatrix} \mathbf{I}_{J_{mob}} & \mathbf{0} & \mathbf{A}_{1,mob} \\ \mathbf{0} & \mathbf{D}_1 \mathbf{A}_{1,het} & \mathbf{D}_1 \tilde{\mathbf{A}}_{1,kin} \end{pmatrix} \begin{pmatrix} \mathbf{r}_{eq} \\ \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix} \quad (19) \\ &= \theta \begin{pmatrix} \mathbf{r}_{mob} + \mathbf{A}_{1,mob} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \\ \mathbf{R} \begin{pmatrix} \mathbf{r}_{sorp} \\ \mathbf{r}_{min} \end{pmatrix} + \mathbf{A}_{1,*} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix}, \end{aligned}$$

where

$$\mathbf{R} := \mathbf{D}_1 \mathbf{A}_{1,het}, \quad \mathbf{A}_{1,*} := \mathbf{D}_1 \tilde{\mathbf{A}}_{1,kin} \quad (20)$$

and where $\mathbf{r}_{eq} = (\mathbf{r}_{mob}, \mathbf{r}_{sorp}, \mathbf{r}_{min})$ analogously to the decomposition of $\mathbf{S}_{1,eq}$ in (8). The size of \mathbf{R} is $N_* \times J_{het}$ and the size of $\mathbf{A}_{1,*}$ is $N_* \times J_{kin}$.

The formulation (19) motivates to introduce the following transformed variables

$$\boldsymbol{\eta} := \tilde{\mathbf{S}}_1^{\perp T} \mathbf{c}, \quad \boldsymbol{\xi} = (\boldsymbol{\xi}_{mob}, \boldsymbol{\xi}_*) := \mathbf{C}_1^T \mathbf{c}. \quad (21)$$

⁴The simplest choice is the following: Let the columns of $\tilde{\mathbf{S}}_1^\perp$ be a basis of the orthogonal complement of the space spanned by the columns of \mathbf{S}_1^* , and let $\mathbf{C}_1 = \mathbf{S}_1^* (\mathbf{S}_1^{*T} \mathbf{S}_1^*)^{-T}$; then $\mathbf{C}_1^T \mathbf{S}_1^* = \mathbf{I}_{J_{mob} + J_{1,het}^* + J_{1,kin}^*}$, i.e., $\mathbf{D}_1 = \mathbf{I}_{N_*}$ ($N_* = J_{1,het}^* + J_{1,kin}^*$). This choice is underlying to the method in [14].

The number of the variables $\boldsymbol{\eta}$ (reaction invariants) is $I - J_{mob} - N_*$, the number of the variables $\boldsymbol{\xi}_{mob}$ is J_{mob} and the number of the variables $\boldsymbol{\xi}_*$ is N_* . Using these new variables (19) reads

$$\begin{aligned}\partial_t(\theta\boldsymbol{\eta}) + L\boldsymbol{\eta} &= \mathbf{0}, \\ \partial_t(\theta\boldsymbol{\xi}_*) + L\boldsymbol{\xi}_* &= \theta\mathbf{R} \begin{pmatrix} \mathbf{r}_{sorp} \\ \mathbf{r}_{min} \end{pmatrix} + \theta\mathbf{A}_{1,*}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}})\end{aligned}\quad (22)$$

and

$$\partial_t(\theta\boldsymbol{\xi}_{mob}) + L\boldsymbol{\xi}_{mob} = \theta\mathbf{r}_{mob} + \theta\mathbf{A}_{1,mob}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}).\quad (23)$$

This last equation is usually dropped together with the unknowns \mathbf{r}_{mob} , which reduces the size of the problem. If desired, this equation can be used for an a posteriori computation of the rates \mathbf{r}_{mob} .

The transformation inverse to (21) obviously reads

$$\mathbf{c} = \begin{pmatrix} \mathbf{C}_1 & \tilde{\mathbf{S}}_1^\perp \end{pmatrix}^{-T} \begin{pmatrix} \boldsymbol{\xi}_{mob} \\ \boldsymbol{\xi}_* \\ \boldsymbol{\eta} \end{pmatrix}.\quad (24)$$

Because of the conditions 2. and 3. the inverse exists. Partitioning the columns of the inverse analogously to the entries of the vector $(\boldsymbol{\xi}_{mob}, \boldsymbol{\xi}_*, \boldsymbol{\eta})$ we write

$$\begin{pmatrix} \mathbf{C}_1 & \tilde{\mathbf{S}}_1^\perp \end{pmatrix}^{-T} = \begin{pmatrix} \mathbf{X} & \mathbf{Y} & \mathbf{Z} \end{pmatrix}.\quad (25)$$

where the size of the blocks $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ is $I \times J_{mob}, I \times N_*, I \times (I - J_{mob} - N_*)$, respectively. With help of the conditions 4. and 5. one gets

$$\begin{pmatrix} \mathbf{C}_1 & \tilde{\mathbf{S}}_1^\perp \end{pmatrix}^T \mathbf{S}_1^* = \begin{pmatrix} \mathbf{C}_1^T \mathbf{S}_1^* \\ \tilde{\mathbf{S}}_1^{\perp T} \mathbf{S}_1^* \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{J_{mob}} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_1 \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Multiplying this equation by (25) from the left we get

$$\mathbf{S}_1^* = \begin{pmatrix} \mathbf{X} & \mathbf{Y} \mathbf{D}_1 \end{pmatrix}.$$

We multiply by \mathbf{A}_1 and obtain (cf. (16))

$$\mathbf{S}_1 = \begin{pmatrix} \mathbf{X} & \mathbf{Y} \mathbf{D}_1 \end{pmatrix} \mathbf{A}_1.$$

Using the block structure of \mathbf{A}_1 (see (17)) and the definition of the matrices \mathbf{R} and $\mathbf{A}_{1,*}$ (see (20)) we get

$$\mathbf{S}_1 = \begin{pmatrix} \mathbf{X} & \mathbf{Y} \mathbf{R} & \mathbf{X} \mathbf{A}_{1,mob} + \mathbf{Y} \mathbf{A}_{1,*} \end{pmatrix}.$$

In particular (see (11) for the block structure of \mathbf{S}_1),

$$\mathbf{X} = \mathbf{S}_{1,mob}, \quad \mathbf{Y}\mathbf{R} = \mathbf{S}_{1,het} = \begin{pmatrix} \mathbf{S}_{1,sorp} & \mathbf{S}_{1,min} \end{pmatrix}. \quad (26)$$

With the left equality and with (25) the inverse transformation (24) can be rewritten as

$$\mathbf{c} = \mathbf{S}_{1,mob}\boldsymbol{\xi}_{mob} + \mathbf{Y}\boldsymbol{\xi}_* + \mathbf{Z}\boldsymbol{\eta}. \quad (27)$$

Transformation of the ODE block

We construct a matrix \mathbf{S}_2^\perp and a matrix \mathbf{B}_2 . The matrix \mathbf{S}_2^\perp consists of a maximal system of linear independent vectors that are orthogonal to all columns of \mathbf{S}_2^* . We choose a matrix \mathbf{B}_2 being of the same size as \mathbf{S}_2^* , that fulfills the condition that the columns of \mathbf{B}_2 , \mathbf{S}_2^\perp form a basis of the whole space. Furthermore \mathbf{B}_2 should have a block structure like \mathbf{S}_2^* (c.f. (15)), i.e.,

$$\mathbf{B}_2 = \begin{pmatrix} * & \mathbf{0} & * \\ \mathbf{0} & \mathbf{I}_{J_{min}} & \mathbf{0} \end{pmatrix}. \quad (28)$$

Note that one possible choice obviously is $\mathbf{B}_2 = \mathbf{S}_2^*$. Analogously to \mathbf{S}_2^\perp a matrix \mathbf{B}_2^\perp is constructed from \mathbf{B}_2 , such that the columns of \mathbf{B}_2^\perp form a maximal system of columns orthogonal to the columns of \mathbf{B}_2 . The orthogonality relations imply

$$\mathbf{S}_2^{\perp T} \mathbf{S}_2^* = \mathbf{0}, \quad \mathbf{B}_2^T \mathbf{B}_2^\perp = \mathbf{0}. \quad (29)$$

Multiplying the block of the ODEs (13) by⁵

$$(\mathbf{S}_2^{\perp T} \mathbf{B}_2^\perp)^{-1} \mathbf{S}_2^{\perp T} \text{ and } (\mathbf{B}_2^T \mathbf{S}_2^*)^{-1} \mathbf{B}_2^T$$

and applying (16), (29) leads to

$$\begin{aligned} \partial_t \left(\theta (\mathbf{S}_2^{\perp T} \mathbf{B}_2^\perp)^{-1} \mathbf{S}_2^{\perp T} \bar{\mathbf{c}} \right) &= \mathbf{0} \\ \partial_t \left(\theta (\mathbf{B}_2^T \mathbf{S}_2^*)^{-1} \mathbf{B}_2^T \bar{\mathbf{c}} \right) &= \theta \mathbf{A}_2 \begin{pmatrix} \mathbf{r}_{eq} \\ \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \end{pmatrix} \end{aligned} \quad (30)$$

This motivates the following definition of new variables:

$$\bar{\boldsymbol{\eta}} := (\mathbf{S}_2^{\perp T} \mathbf{B}_2^\perp)^{-1} \mathbf{S}_2^{\perp T} \bar{\mathbf{c}}, \quad \bar{\boldsymbol{\xi}} := (\mathbf{B}_2^T \mathbf{S}_2^*)^{-1} \mathbf{B}_2^T \bar{\mathbf{c}} \quad (31)$$

⁵For a proof that the two inverses exist see [9, Sec. 3.1].

The variables $\bar{\xi}$ are partitioned analogously to the partitioning of the columns of \mathbf{S}_2^* in (15) into

$$\bar{\xi} = (\bar{\xi}_{sorp}, \bar{\xi}_{min}, \bar{\xi}_{kin}) =: (\bar{\xi}_{het}, \bar{\xi}_{kin}). \quad (32)$$

Then (30) reads

$$\partial_t(\theta\bar{\eta}) = \mathbf{0} \quad (33)$$

$$\partial_t(\theta\bar{\xi}_{sorp}) = \theta(\mathbf{r}_{sorp} + \mathbf{A}_{2,sorp}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}})) \quad (34)$$

$$\partial_t(\theta\bar{\xi}_{min}) = \theta\mathbf{r}_{min} \quad (35)$$

$$\partial_t(\theta\bar{\xi}_{kin}) = \theta\mathbf{A}_{2,kin}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}). \quad (36)$$

Hence, we have found some more reaction invariants (33). The variables $\bar{\eta}, \bar{\xi}_{sorp}, \bar{\xi}_{min}, \bar{\xi}_{kin}$ have the sizes $\bar{I} - J_{sorp} - J_{kin} - J_{2,kin}^*, J_{sorp}, J_{min}, J_{2,kin}^*$.

For the transformation inverse to (31) we get

$$\bar{\mathbf{c}} = \mathbf{S}_2^*\bar{\xi} + \mathbf{B}_2^\perp\bar{\eta}. \quad (37)$$

To verify that (37) and (31) are inverse, just substitute the one into the other and use the orthogonality relations (29).

Since the matrices \mathbf{S}_2^* and \mathbf{B}_2 are of the form (28), the transformed variables $\bar{\xi}_{sorp}$ and $\bar{\xi}_{kin}$ in (31) only depend on the upper part $\bar{\mathbf{c}}_{nmin}$ of $\bar{\mathbf{c}}$, but not on the lower part $\bar{\mathbf{c}}_{min}$. Furthermore, since the second column block of (28) consists of unit vectors, the orthogonality relations (29) lead to the fact that the last J_{min} entries in every column of \mathbf{S}_2^\perp and \mathbf{B}_2^\perp are always zero. Hence, also $\bar{\eta}$ in (31) only depends on $\bar{\mathbf{c}}_{nmin}$, but not on $\bar{\mathbf{c}}_{min}$, and we can decompose \mathbf{B}_2^\perp into an upper block $\tilde{\mathbf{B}}_2^\perp$ and a lower zero block. Using this and the partitionings (15), (32) we can rewrite the inverse transformation (37) as

$$\bar{\mathbf{c}} = \begin{pmatrix} \mathbf{S}_{2,sorp}\bar{\xi}_{sorp} + \mathbf{S}_{2,kin}^*\bar{\xi}_{kin} + \tilde{\mathbf{B}}_2^\perp\bar{\eta} \\ \bar{\xi}_{min} \end{pmatrix}. \quad (38)$$

By now, we have transformed system (12)-(14) into (22), (33)-(36), (14), for the unknowns $\boldsymbol{\eta}, \bar{\boldsymbol{\eta}}, \boldsymbol{\xi}_{mob}, \boldsymbol{\xi}_*, \bar{\boldsymbol{\xi}}_{het} = (\bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}), \bar{\boldsymbol{\xi}}_{kin}$, where \mathbf{c} and $\bar{\mathbf{c}}$ are expressed by (24)/(27), (38). The equations for $\boldsymbol{\eta}$ and for $\bar{\boldsymbol{\eta}}$ are decoupled and linear, reducing the size of the remaining nonlinear system. Note that, while the equilibrium rates \mathbf{r}_{mob} are already eliminated (cf. (23)), our aim will be to eliminate the equilibrium rates $\mathbf{r}_{sorp}, \mathbf{r}_{min}$ as well.

Additional variables

Before we present the final transformed system of equations, let us introduce additional variables $\tilde{\xi}$ defined by

$$\tilde{\xi} := \xi_* - R\bar{\xi}_{het}. \quad (39)$$

The advantage of these additional variables is discussed in Sec. 4.5.

If we solve equation (39) for ξ_* and substitute it into the inverse transformation (27), then using (26), the inverse transformation reads

$$\mathbf{c} = \mathbf{S}_{1,mob}\xi_{mob} + \mathbf{Y}\tilde{\xi} + \mathbf{S}_{1,het}\bar{\xi}_{het} + \mathbf{Z}\eta. \quad (40)$$

Solving the ODEs (34) and (35) for \mathbf{r}_{sorp} and \mathbf{r}_{min} and plugging this in the PDEs (22) gives

$$\begin{aligned} & \partial_t(\theta\xi_*) + L\xi_* \\ &= \mathbf{R} \begin{pmatrix} \partial_t(\theta\bar{\xi}_{sorp}) - \theta\mathbf{A}_{2,sorp}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \\ \partial_t(\theta\bar{\xi}_{min}) \end{pmatrix} + \theta\mathbf{A}_{1,*}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}), \end{aligned}$$

which can be rewritten as

$$\partial_t(\theta\tilde{\xi}) + L\xi_* = \theta\mathbf{A}_*\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \quad (41)$$

with

$$\mathbf{A}_* := \mathbf{A}_{1,*} - \mathbf{R} \begin{pmatrix} \mathbf{A}_{2,sorp} \\ \mathbf{0} \end{pmatrix}. \quad (42)$$

Note that all the equilibrium rates are eliminated now from the system. Altogether we obtain the equations

$$\partial_t(\theta\eta) + L\eta = \mathbf{0} \quad (43)$$

$$\partial_t(\theta\bar{\eta}) = \mathbf{0} \quad (44)$$

$$\tilde{\xi} = \xi_* - R\bar{\xi}_{het} \quad (45)$$

$$\partial_t(\theta\tilde{\xi}) + L\xi_* = \theta\mathbf{A}_*\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \quad (46)$$

$$\partial_t(\theta\bar{\xi}_{kin}) = \theta\mathbf{A}_{2,kin}\mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) \quad (47)$$

$$\phi(\mathbf{c}, \bar{\mathbf{c}}) = \mathbf{0} \quad (48)$$

with

$$\begin{aligned} \mathbf{c} &= \mathbf{S}_{1,mob}\xi_{mob} + \mathbf{Y}\tilde{\xi} + \mathbf{S}_{1,het}\bar{\xi}_{het} + \mathbf{Z}\eta \\ \bar{\mathbf{c}} &= \begin{pmatrix} \mathbf{S}_{2,sorp}\bar{\xi}_{sorp} + \mathbf{S}_{2,kin}^*\bar{\xi}_{kin} + \tilde{\mathbf{B}}_2^\perp\bar{\eta} \\ \bar{\xi}_{min} \end{pmatrix}. \end{aligned} \quad (49)$$

Compared to (12)-(13), but also to other reformulations used in the literature, the advantage of system (43)-(48) is that the equations (43) and (44) decouple from the rest of the system. So the coupled nonlinear system, which must be solved numerically, is smaller. Counting the PDEs in the nonlinear system (the number of PDEs severely influences the number of nonzero entries in the Jacobian), we find N_* many. As pointed out in Footnote 4, a proper choice of \mathbf{S}_1^\perp leads to $N_* = J_{1,het}^* + J_{1,kin}^* \leq J_{sorp} + J_{min} + J_{kin} = J - J_{mob}$.

4 Implicit elimination

4.1 The main idea

In order to further reduce the size of the system, we are formally solving some 'local' equations (equations without spatial couplings) for some 'local' variables, and substitute these local variables into the remaining 'global' equations (PDEs). Hence, we perceive the system (45)-(48) as

$$\begin{aligned}\mathbf{f}(\mathbf{u}, \mathbf{v}) &= \mathbf{0} \\ \mathbf{g}(\mathbf{u}, \mathbf{v}) &= \mathbf{0}\end{aligned}\tag{50}$$

with global equations \mathbf{f} , local equations \mathbf{g} , global variables \mathbf{u} , local variables \mathbf{v} . The number of global (local) variables must equal the number of global (local) equations.

If the second block has a resolution function $\mathbf{v} = \mathbf{v}(\mathbf{u})$, i.e., the function $\mathbf{v}(\mathbf{u})$ fulfills

$$\mathbf{g}(\mathbf{u}, \mathbf{v}(\mathbf{u})) = \mathbf{0} \text{ for all } \mathbf{u},\tag{51}$$

then we can rewrite system (50) by the smaller system

$$\mathbf{f}(\mathbf{u}, \mathbf{v}(\mathbf{u})) = \mathbf{0}.\tag{52}$$

The existence of the resolution function $\mathbf{v}(\mathbf{u})$ will be checked in Sec. 4.2. To solve this condensed system (52) for \mathbf{u} with Newton's method, the evaluation of the mapping $\mathbf{u} \mapsto \mathbf{f}(\mathbf{u}, \mathbf{v}(\mathbf{u}))$, i.e., of $\mathbf{v}(\mathbf{u})$, and of the Jacobian of (52), i.e.,⁶

$$\mathbf{J} = \partial_1 \mathbf{f} + \partial_2 \mathbf{f} \mathbf{v}',\tag{53}$$

are needed for a given Newton iterate \mathbf{u} . The evaluation of $\mathbf{v}(\mathbf{u})$ can be done by solving the local problem

$$\mathbf{g}(\mathbf{u}, \mathbf{v}) = \mathbf{0} \text{ for given } \mathbf{u}$$

⁶ $\partial_{1/2}$ denote the partial derivative with respect to the first/second block of variables

with Newton's method, i.e., by a nested Newton iteration. Note that due to the lack of spatial couplings within \mathbf{g} , there is a decoupling in the sense that the nested Newton iteration can be done mesh point by mesh point. The derivative \mathbf{v}' is obtained by differentiating (51) with respect to \mathbf{u} . This yields $\partial_1 \mathbf{g} + \partial_2 \mathbf{g} \mathbf{v}' = \mathbf{0}$, i.e.,

$$\mathbf{v}' = -(\partial_2 \mathbf{g})^{-1} \partial_1 \mathbf{g}. \quad (54)$$

See also Sec. 4.4 for some details on the computation of (54). The invertibility of matrix $\partial_2 \mathbf{g}$ is directly related to the existence of the resolution function, cf. Sec. 4.2. Note that $\partial_2 \mathbf{g}$ is also the Jacobian of the local problem.

4.2 Existence of the resolution function

We consider equations (47)-(48) to be local and (45)-(46) to be global in the sense of Sec. 4.1. The unknowns are split into local unknowns

$$\mathbf{v} = \boldsymbol{\xi}_{loc} := (\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \bar{\boldsymbol{\xi}}_{kin}) \quad (55)$$

and global unknowns

$$\mathbf{u} = \boldsymbol{\xi}_{glob} := (\tilde{\boldsymbol{\xi}}, \boldsymbol{\xi}_*). \quad (56)$$

Now we want to prove the existence of a resolution function $\boldsymbol{\xi}_{loc}(\boldsymbol{\xi}_{glob})$ defined by the equations (47)-(48) (after time discretization of (47) with the implicit Euler method) and $(\mathbf{c}, \bar{\mathbf{c}})$ given by (49). Note that the resolution function $\boldsymbol{\xi}_{loc}(\boldsymbol{\xi}_{glob})$ does not depend on the variables $\boldsymbol{\xi}_*$ because these variables do not appear in the equations (47)-(48) and in the inverse transformation (49). Hence, we want to show that the equations (48) and

$$\phi_{kin} := \bar{\boldsymbol{\xi}}_{kin} - \bar{\boldsymbol{\xi}}_{kin,old} - \Delta t \mathbf{A}_{2,kin} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}) = 0,$$

which is the time-discretized version of (47), have a resolution function

$$\tilde{\boldsymbol{\xi}} \mapsto (\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \bar{\boldsymbol{\xi}}_{kin})$$

We assume that a positive lower bound for the concentrations exists. The following proof is adapted from [14, Appendix]. To apply the implicit function theorem, we have to check that the matrix

$$\frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min}, \phi_{kin})}{\partial(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \bar{\boldsymbol{\xi}}_{kin})}, \quad (57)$$

which is $\partial_2 \mathbf{g}$, is invertible. For that purpose the equilibrium mineral reactions are split into inactive (index \mathcal{I}) and active (index \mathcal{A}) reactions:

$$\mathbf{S}_{1,min} = \left(\mathbf{S}_{1,min,\mathcal{I}} \quad \mathbf{S}_{1,min,\mathcal{A}} \right)$$

A reaction is called inactive when in $\min\{-\ln(K_j) + \sum_{i=1}^I s_{min,ij} \ln(c_i), \bar{c}_{min,j}\}$ the minimum is attained in the first argument and otherwise called active.

In a first step, we consider matrix (57) in the limit case $\Delta t = 0$. For this case we obtain that the matrix (57) reads as in Table 1, where we have used the diagonal matrices

$$\begin{aligned} \mathbf{\Lambda} &= \text{diag}(1/c_1, \dots, 1/c_I), \\ \bar{\mathbf{\Lambda}}_{nmin} &= \text{diag}(1/\bar{c}_{I+1}, \dots, 1/\bar{c}_{I+\bar{I}_{nmin}}). \end{aligned}$$

We can rewrite that matrix from Table 1 as

Table 1: Jacobian (57) of the local problem in the limit case $\Delta t = 0$.

$$\begin{aligned} & \frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min}, \phi_{kin})}{\partial(\xi_{mob}, \bar{\xi}_{sorp}, \bar{\xi}_{min}, \bar{\xi}_{kin})} \\ &= \begin{pmatrix} \frac{\partial \phi_{mob}}{\partial \mathbf{c}} \mathbf{S}_{1,mob} & \frac{\partial \phi_{mob}}{\partial \mathbf{c}} \mathbf{S}_{1,sorp} & \frac{\partial \phi_{mob}}{\partial \mathbf{c}} \mathbf{S}_{1,min} & \mathbf{0} \\ \frac{\partial \phi_{sorp}}{\partial \mathbf{c}} \mathbf{S}_{1,mob} & \frac{\partial \phi_{sorp}}{\partial \mathbf{c}} \mathbf{S}_{1,sorp} + \frac{\partial \phi_{sorp}}{\partial \bar{\mathbf{c}}_{nmin}} \mathbf{S}_{2,sorp} & \frac{\partial \phi_{sorp}}{\partial \mathbf{c}} \mathbf{S}_{1,min} & \mathbf{0} \\ \frac{\partial \phi_{min}}{\partial \mathbf{c}} \mathbf{S}_{1,mob} & \frac{\partial \phi_{min}}{\partial \mathbf{c}} \mathbf{S}_{1,sorp} & \frac{\partial \phi_{min}}{\partial \mathbf{c}} \mathbf{S}_{1,min} + \frac{\partial \phi_{min}}{\partial \bar{\mathbf{c}}_{min}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{S}_{1,mob}^T \mathbf{\Lambda} \mathbf{S}_{1,mob} & \mathbf{S}_{1,mob}^T \mathbf{\Lambda} \mathbf{S}_{1,sorp} & \mathbf{S}_{1,mob}^T \mathbf{\Lambda} \mathbf{S}_{1,min} & \mathbf{0} \\ \mathbf{S}_{1,sorp}^T \mathbf{\Lambda} \mathbf{S}_{1,mob} & \mathbf{S}_{1,sorp}^T \mathbf{\Lambda} \mathbf{S}_{1,sorp} + \mathbf{S}_{2,sorp}^T \bar{\mathbf{\Lambda}}_{nmin} \mathbf{S}_{2,sorp} & \mathbf{S}_{1,sorp}^T \mathbf{\Lambda} \mathbf{S}_{1,min} & \mathbf{0} \\ \mathbf{S}_{1,min,\mathcal{I}}^T \mathbf{\Lambda} \mathbf{S}_{1,mob} & \mathbf{S}_{1,min,\mathcal{I}}^T \mathbf{\Lambda} \mathbf{S}_{1,sorp} & \mathbf{S}_{min,\mathcal{I}}^T \mathbf{\Lambda} \mathbf{S}_{1,min} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \begin{pmatrix} \mathbf{0} & \mathbf{I}_{\mathcal{A}} \end{pmatrix} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} \end{pmatrix} \end{aligned}$$

$$\frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min}, \phi_{kin})}{\partial(\xi_{mob}, \bar{\xi}_{sorp}, \bar{\xi}_{min}, \bar{\xi}_{kin})} = \begin{pmatrix} \mathbf{M}_1 & \mathbf{M}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{\mathcal{A}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{J_{2,kin}^*} \end{pmatrix} \quad (58)$$

with

$$\begin{aligned} \mathbf{M}_1 &= \mathbf{Q}^T \tilde{\mathbf{\Lambda}} \mathbf{Q}, \quad \tilde{\mathbf{\Lambda}} = \begin{pmatrix} \mathbf{\Lambda} & \mathbf{0} \\ \mathbf{0} & \bar{\mathbf{\Lambda}}_{nmin} \end{pmatrix}, \\ \mathbf{Q} &= \begin{pmatrix} \mathbf{S}_{1,mob} & \mathbf{S}_{1,sorp} & \mathbf{S}_{1,min,\mathcal{I}} \\ \mathbf{0} & \mathbf{S}_{2,sorp} & \mathbf{0} \end{pmatrix}, \end{aligned} \quad (59)$$

$$\mathbf{M}_2 = \begin{pmatrix} \mathbf{S}_{1,mob}^T \\ \mathbf{S}_{1,sorp}^T \\ \mathbf{S}_{1,min,\mathcal{I}}^T \end{pmatrix} \mathbf{\Lambda} \mathbf{S}_{1,min,\mathcal{A}}.$$

Due to the linear independence assumption (see after (10)) the columns of the matrix \mathbf{Q} are linear independent. Hence, matrix \mathbf{M}_1 is symmetric positive definite because the diagonal matrix $\tilde{\mathbf{\Lambda}}_{nmin}$ only has positive entries. Hence, matrix (58) is invertible.

For sufficiently small Δt we use a continuity argument, exactly as in [14], to finish the proof.

Using this resolution function and assuming that the decoupled equations (43) and (44) are solved, (43)-(48) can be reduced to the condensed system (cf. (52))

$$\tilde{\boldsymbol{\xi}} = \boldsymbol{\xi}_* - \mathbf{R}\bar{\boldsymbol{\xi}}_{het}(\tilde{\boldsymbol{\xi}}) \quad (60)$$

$$\partial_t(\theta\tilde{\boldsymbol{\xi}}) + L\tilde{\boldsymbol{\xi}}_* = \theta\mathbf{A}_*\tilde{\mathbf{r}}_{kin}(\tilde{\boldsymbol{\xi}}). \quad (61)$$

The reaction rates \mathbf{r}_{kin} are written as a function of $\tilde{\boldsymbol{\xi}}$. This can be achieved by (49) and the resolution function. Comparing with (43)-(48) one can see that the number of coupled equations is reduced drastically. This nonlinear system will be used for the numerical computations. Contrary to DSA there is no nonlinearity under the transport operator. That is an advantage because it is known that a nonlinearity under the differential operator may cause numerical problems; at least it produces many nonzero entries in the Jacobian.

4.3 The local problem seen as a minimization problem

It is well known that mass action reactions are closely related to the minimization of the so-called Gibbs free energy. We want to show that in our formulation the local problem can be regarded as a convex minimization problem, provided that there are no variables $\tilde{\boldsymbol{\xi}}_{kin}$ (i.e., $J_{2,kin}^* = 0$). In [11, Sec. 2.4.4] this consideration is done for the special case that there are no equilibrium minerals and for a formulation that comes without the additional variables $\tilde{\boldsymbol{\xi}}$.

First we define the functional

$$G(\mathbf{c}, \bar{\mathbf{c}}) := \sum_{i=1}^I \mu_i(c_i)c_i + \sum_{i=I+1}^{I+\bar{I}} \bar{\mu}_i(\bar{c}_i)\bar{c}_i$$

where

$$\begin{aligned}\mu_i(c_i) &:= \mu_{0,i} - 1 + \ln(c_i) \\ \bar{\mu}_i(\bar{c}_i) &:= \begin{cases} \bar{\mu}_{0,i} - 1 + \ln(\bar{c}_{nmin,i}), & i = I+1, \dots, I+\bar{I}_{nmin} \\ \bar{\mu}_{0,i}, & i = I+\bar{I}_{nmin}+1, \dots, I+\bar{I}, \end{cases}\end{aligned}$$

where the vector $(\boldsymbol{\mu}_0, \bar{\boldsymbol{\mu}}_0) \in \mathbb{R}^{I+\bar{I}}$ is a solution of the linear system

$$\mathbf{S}_{eq}^T \begin{pmatrix} \boldsymbol{\mu}_0 \\ \bar{\boldsymbol{\mu}}_0 \end{pmatrix} = -\ln(\mathbf{K}). \quad (62)$$

We calculate

$$\nabla G(\mathbf{c}, \bar{\mathbf{c}}) = (\boldsymbol{\mu}_0, \bar{\boldsymbol{\mu}}_{0,nmin}, \bar{\boldsymbol{\mu}}_{0,min}) + (\ln(\mathbf{c}), \ln(\bar{\mathbf{c}}_{nmin}), \mathbf{0}).$$

We see that the Hessian $D^2G = \text{diag}(\frac{1}{c_1}, \dots, \frac{1}{\bar{c}_{I+\bar{I}_{nmin}}}, 0, \dots, 0)$ is positive semidefinite. So it follows that G as a function of $(\mathbf{c}, \bar{\mathbf{c}})$ is a convex functional.

With help of the inverse transformation (49) and for given values for the variables $\boldsymbol{\eta}, \tilde{\boldsymbol{\xi}}, \bar{\boldsymbol{\eta}}$ we can write G as a function of the local variables $\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}$. Remember that in this subsection we assume that there are no variables $\bar{\boldsymbol{\xi}}_{kin}$. Also as a function of $(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})$ G is a convex function. This can be seen as follows. Using the fact that $\frac{\partial(\mathbf{c}, \bar{\mathbf{c}})}{\partial(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})} = \mathbf{S}_{eq}$, which follows immediately from (49), we compute

$$\begin{aligned}\nabla_{(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})} G &= \mathbf{S}_{eq}^T \nabla_{(\mathbf{c}, \bar{\mathbf{c}})} G \\ &= \mathbf{S}_{eq}^T \begin{pmatrix} \boldsymbol{\mu}_0 \\ \bar{\boldsymbol{\mu}}_{0,nmin} \\ \bar{\boldsymbol{\mu}}_{0,min} \end{pmatrix} + \mathbf{S}_{eq}^T \begin{pmatrix} \ln(\mathbf{c}) \\ \ln(\bar{\mathbf{c}}_{nmin}) \\ \mathbf{0} \end{pmatrix}. \end{aligned} \quad (63)$$

For the Hessian of G we get

$$D^2_{(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})} G = \mathbf{S}_{eq}^T \boldsymbol{\Lambda} \mathbf{S}_{eq} \quad (64)$$

with $\boldsymbol{\Lambda} = \text{diag}(\frac{1}{c_1}, \dots, \frac{1}{c_I}, \frac{1}{\bar{c}_{I+1}}, \dots, \frac{1}{\bar{c}_{I+\bar{I}_{nmin}}}, 0, \dots, 0)$. We know that the matrix $\mathbf{S}_{eq}^T \boldsymbol{\Lambda} \mathbf{S}_{eq}$ is positive semidefinite because $\boldsymbol{\Lambda}$ is a diagonal matrix with nonnegative entries. Hence G as a function of $(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})$ is a convex functional.

Now we consider the constrained minimization problem

$$\begin{aligned}\min & G(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}) \\ \text{s.t.} & \bar{\boldsymbol{\xi}}_{min} \geq \mathbf{0}. \end{aligned} \quad (65)$$

The Lagrangian of this minimization problem reads

$$\mathcal{L}(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \boldsymbol{\nu}) = G(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}) - \bar{\boldsymbol{\xi}}_{min} \cdot \boldsymbol{\nu}.$$

Using (63) and (62) we get for the associated KKT system

$$\begin{aligned} 0 &= \nabla_{(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})} \mathcal{L} \\ &= \mathbf{S}_{eq}^T \begin{pmatrix} \boldsymbol{\mu}_0 \\ \bar{\boldsymbol{\mu}}_{0,min} \\ \bar{\boldsymbol{\mu}}_{0,min} \end{pmatrix} + \mathbf{S}_{eq}^T \begin{pmatrix} \ln(\mathbf{c}) \\ \ln(\bar{\mathbf{c}}_{nmin}) \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\boldsymbol{\nu} \end{pmatrix} \\ &= -\ln(\mathbf{K}) + \mathbf{S}_{eq}^T \begin{pmatrix} \ln(\mathbf{c}) \\ \ln(\bar{\mathbf{c}}_{nmin}) \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\boldsymbol{\nu} \end{pmatrix} \\ &= \begin{pmatrix} \phi_{mob}(\mathbf{c}) \\ \phi_{sorp}(\mathbf{c}, \bar{\mathbf{c}}_{nmin}) \\ \psi_{min}(\mathbf{c}) \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ -\boldsymbol{\nu} \end{pmatrix}, \\ &\quad \nu_j \bar{\xi}_{min,j} = 0, \quad \bar{\xi}_{min,j} \geq 0, \quad \nu_j \geq 0, \quad j = 1, \dots, J_{min} \end{aligned}$$

where the entries of $\boldsymbol{\psi}_{min}$ are defined according to (3). Hence the minimization problem (65) is equivalent to solving the equilibrium conditions (48). So the resolution function $\boldsymbol{\xi}_{loc}(\boldsymbol{\xi}_{glob})$ of the previous subsection can also be defined as the solution of the constrained minimization problem (65) if there are no variables $\bar{\boldsymbol{\xi}}_{kin}$. The difference to [11, Sec. 2.4.4], where the additional variables $\bar{\boldsymbol{\xi}}$ are not used, is that we can write the local problem as *one* minimization problem while in [11] there are two coupled minimization problems are needed.

Note that the Hessian (64) can be written as

$$\begin{aligned} D^2_{(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min})} G &= \hat{\mathbf{S}}_{eq}^T \text{diag}\left(\frac{1}{c_1}, \dots, \frac{1}{\bar{c}_{I+\bar{I}_{nmin}}}\right) \hat{\mathbf{S}}_{eq}, \\ \hat{\mathbf{S}}_{eq} &= \begin{pmatrix} \mathbf{S}_{1,mob} & \mathbf{S}_{1,sorp} & \mathbf{S}_{1,min} \\ \mathbf{0} & \mathbf{S}_{2,sorp} & \mathbf{0} \end{pmatrix}. \end{aligned}$$

From the assumption (10) we know that the columns of $\hat{\mathbf{S}}_{eq}$ are linear independent. As a consequence, the Hessian is strictly positive definit, i.e., the constrained minimization problem is strictly convex. From this we can derive that the solution of the local problem is unique.

4.4 Calculation of matrix \mathbf{v}' for the reduction scheme

In order to use the reformulation of Sec. 3 with the implicit elimination of Sec. 4.1-4.2, we have to compute the matrix \mathbf{v}' (see (54)) to apply a Newton

step for (52). Instead of computing \mathbf{v}' exactly, we will use an approximation. Note that the usage of an approximate \mathbf{v}' does not affect the accuracy of the solution of the time step; it might only affect the convergence of the Newton iteration, if the approximation is poor. We assume that kinetic rate terms, since they contain a factor Δt , are sufficiently small to be neglected for the computation of \mathbf{v}' . This means that within the computation of \mathbf{v}' , we replace any occurrence of Δt by zero.

To compute \mathbf{v}' , i.e., $D_{\tilde{\boldsymbol{\xi}}}(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \bar{\boldsymbol{\xi}}_{kin})$, by formula (54) we need the matrices

$$\begin{aligned}\partial_2 \mathbf{g} &= \frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min}, \phi_{kin})}{\partial(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min}, \bar{\boldsymbol{\xi}}_{kin})}, \\ \partial_1 \mathbf{g} &= \frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min}, \phi_{kin})}{\partial \tilde{\boldsymbol{\xi}}}.\end{aligned}$$

With the approximation described above, we obtain $\partial_2 \mathbf{g}$ by the formulas (58)-(59).

Analogously we can calculate that

$$\begin{aligned}-\partial_1 \mathbf{g} &= -\frac{\partial(\phi_{mob}, \phi_{sorp}, \phi_{min, \mathcal{I}}, \phi_{min, \mathcal{A}}, \phi_{kin})}{\partial \tilde{\boldsymbol{\xi}}} \\ &= \begin{pmatrix} \mathbf{Q}^T \tilde{\boldsymbol{\Lambda}} \mathbf{C} \\ \mathbf{0} \\ \mathbf{I}_{J_{2,kin}^*} \end{pmatrix} \quad \text{with } \mathbf{C} = \begin{pmatrix} -\mathbf{Y} \\ \mathbf{0} \end{pmatrix}.\end{aligned}$$

Because of the block structure of $\partial_1 \mathbf{g}$ and of $\partial_2 \mathbf{g}$ we obtain that the matrix \mathbf{v}' has a block structure $\begin{pmatrix} \mathbf{U} \\ \mathbf{0} \\ \mathbf{I}_{J_{2,kin}^*} \end{pmatrix}$ where \mathbf{U} is the solution of the linear systems

$$\mathbf{Q}^T \tilde{\boldsymbol{\Lambda}} \mathbf{Q} \mathbf{U} = \mathbf{Q}^T \tilde{\boldsymbol{\Lambda}} \mathbf{C}, \quad (66)$$

to be solved at every mesh point. \mathbf{U} contains the derivatives $D_{\tilde{\boldsymbol{\xi}}}(\boldsymbol{\xi}_{mob}, \bar{\boldsymbol{\xi}}_{sorp}, \bar{\boldsymbol{\xi}}_{min, \mathcal{I}})$. The matrix of the linear systems $\mathbf{Q}^T \tilde{\boldsymbol{\Lambda}} \mathbf{Q}$ is symmetric and positive definite and hence it is always invertible.

4.5 Benefit of the Additional Variables

In the end of Sec. 3 we introduced the additional variables $\tilde{\boldsymbol{\xi}}$, which were not present in earlier versions of the reformulation method. Let us motivate this.

We need the matrix $\mathbf{v}' = D_{\xi_{glob}} \xi_{loc}$ to solve the condensed system (60)-(61). With introduction of the additional variables, it turns out that the entries of this matrix, i.e., the entries of \mathbf{U} (cf. Sec. 4.4) are bounded independently of concentration values and reaction constants (see [9, Sec. 3.5.1] for a proof). This would not be the case for matrix \mathbf{v}' without the additional variables ξ . As a consequence of the boundedness, for the MoMaS model problem, one equilibrium sorption reaction and in the limit case $\Delta t = 0$, one can prove for the condition number of the Jacobian matrix (53) that $\kappa(\mathbf{J}) \leq (1 + \sqrt{5})^2/4$ (see [9, Sec. 3.5.2]). Without the additional variables $\tilde{\xi}$ the condition number for large K (K : nonlogarithmized reaction constant) is $\kappa(\mathbf{J}) \approx 4K$ (see [9, Sec. 3.6.1]). In realistic problems reaction constants can have very large values, such as $K = 10^{35}$. So the additional variables are mandatory for many realistic problems.

5 Choice of the Transformation Matrices

In the following we will give two specific choices for the transformation matrices $\mathbf{C}_1, \tilde{\mathbf{S}}_1^\perp, \mathbf{B}_2, \mathbf{B}_2^\perp$ of Sec. 3.

The first choice allows to establish a close connection to the well known Morel formulation. Hence, the computer code for the reformulation technique can be used to simulate Morel-based algorithms.

The second method minimizes the number of remaining coupled PDEs and coincides with the reformulation presented in [9] and is similar to the version given in [14].

5.1 Morel formulation as special case of general reduction scheme

If there are no kinetic reactions, then in the Morel formulation the stoichiometric matrix has the form (see [10] for the notation used in the Morel formulation)

$$\mathbf{S} = \begin{pmatrix} \mathbf{C} & \mathbf{A} & \mathbf{D} \\ -\mathbf{I}_{J_{mob}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{B}} & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_{J_{sorp}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & -\mathbf{I}_{J_{min}} \end{pmatrix},$$

and some variables $\mathbf{T}, \mathbf{T}_M, \mathbf{T}_F, \mathbf{W}$ denoting the so-called total, total mobile, total fixed and constant concentrations. It is possible to achieve that

$$\tilde{\xi} = \mathbf{T}, \quad \xi_* = \mathbf{T}_M, \quad \mathbf{R}\bar{\xi}_{het} = -\mathbf{T}_F, \quad \bar{\eta} = \mathbf{W},$$

giving our variables $\tilde{\xi}, \xi_*, \bar{\eta}$ the concrete meaning of the total, the total mobile, and the constant concentrations.

For this purpose one has to choose the transformation matrices

$$\mathbf{C}_1 = \begin{pmatrix} \mathbf{0} & \mathbf{I}_{N_*} \\ -\mathbf{I}_{J_{mob}} & \mathbf{C}^T \end{pmatrix}, \quad \tilde{\mathbf{S}}_1^\perp \text{ empty matrix} \quad (67)$$

$$\mathbf{B}_2 = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{I}_{J_{sorp}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{J_{min}} \end{pmatrix}, \quad \mathbf{B}_2^\perp = \begin{pmatrix} \mathbf{I}_{\bar{I}-J_{sorp}-J_{min}} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad (68)$$

with $N_* = I - J_{mob}$. Note that the matrix $\tilde{\mathbf{S}}_1^\perp$ consists of zero columns, and so there will be no decoupled PDEs. This is of course a disadvantage for efficient numerical computations. For the variables ξ defined in (21) it holds

$$\begin{aligned} \begin{pmatrix} \xi_{mob} \\ \xi_* \end{pmatrix} &= \mathbf{C}_1^T \mathbf{c} = \begin{pmatrix} \mathbf{0} & -\mathbf{I}_{J_{mob}} \\ \mathbf{I}_{N_*} & \mathbf{C} \end{pmatrix} \begin{pmatrix} \mathbf{c}_{prim} \\ \mathbf{c}_{sec} \end{pmatrix} \\ &= \begin{pmatrix} -\mathbf{c}_{sec} \\ \mathbf{c}_{prim} + \mathbf{C}\mathbf{c}_{sec} \end{pmatrix}. \end{aligned}$$

In particular, one gets the coincidence $\xi_* = \mathbf{c}_{prim} + \mathbf{C}\mathbf{c}_{sec} = \mathbf{T}_M$. The above choice of $\mathbf{B}_2, \mathbf{B}_2^\perp$ leads to

$$\bar{\xi}_{sorp} = -\bar{c}_{nmin,sec}, \quad \bar{\xi}_{min} = -\bar{c}_{min}, \quad \bar{\eta} = \mathbf{W}.$$

Since all columns of \mathbf{S}_1 are linear independent there is $\mathbf{S}_1^* = \mathbf{S}_1$ and $\mathbf{A}_1 = \mathbf{I}_{J_{eq}}$. So by computing the matrix product

$$\begin{aligned} \mathbf{C}_1^T \mathbf{S}_1^* &= \begin{pmatrix} \mathbf{0} & -\mathbf{I}_{J_{mob}} \\ \mathbf{I}_{N_*} & \mathbf{C} \end{pmatrix} \begin{pmatrix} \mathbf{C} & \mathbf{A} & \mathbf{D} \\ -\mathbf{I}_{J_{mob}} & \mathbf{0} & \mathbf{0} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{I}_{J_{mob}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} & \mathbf{D} \end{pmatrix} \end{aligned}$$

ones sees with the definitions of the matrices \mathbf{D}_1 (18) and \mathbf{R} (20) that

$$\mathbf{R} = \begin{pmatrix} \mathbf{A} & \mathbf{D} \end{pmatrix}.$$

Using this we obtain for $\tilde{\xi}$ defined in (39) that

$$\begin{aligned}\tilde{\xi} &= \xi_* - R\bar{\xi}_{het} = \xi_* - (A \ D) \bar{\xi}_{het} \\ &= c_{prim} + Cc_{sec} + A\bar{c}_{nmin,sec} + D\bar{c}_{min}.\end{aligned}$$

Therefore $\tilde{\xi} = T$ and $R\bar{\xi}_{het} = -T_F$ hold.

With help of these identities it is easy to see that also the Morel equations

$$\begin{aligned}\partial_t(\theta T) + LT_M &= \mathbf{0} \\ T &= T_M + T_F(\bar{c}_{nmin,sec}, \bar{c}_{min}) \\ \partial_t(\theta W) &= \mathbf{0} \\ \phi(c, \bar{c}) &= \mathbf{0}\end{aligned}$$

coincide with ours. Using the reduction scheme one has a resolution function $\bar{\xi}_{het}(\tilde{\xi})$. Using the identities between the total concentrations and the transformed variables one gets a resolution function $T_F(T)$.

Plugging in this resolution function gives

$$\begin{aligned}\partial_t(\theta T) + LT_M &= \mathbf{0} \\ T &= T_M + T_F(T).\end{aligned}$$

Solving this system is called *global-ODE approach* (see [8]). The formulation used here is very similar to the formulation in [1]. The only difference is that the equation $T_F = \psi(T)$, appearing in the formulation of [1], is plugged into the other equations of the formulation.

5.2 The standard version of the reduction scheme as special case of the general reduction scheme

To get the standard version of the reduction scheme, as it is presented in [9, 10], and similar in [13, 11], proceed as follows. Choose a matrix S_1^\perp such that its columns form a basis of the orthogonal complement of the range of S_1^* . Then let B_1 be a matrix such that the columns of B_1 and of S_1^\perp form a basis of \mathbb{R}^I (so $B_1 = S_1^*$ is a possible choice)⁷. Then construct B_1^\perp from B_1 in the same way S_1^\perp is constructed from S_1^* . Then set

$$\tilde{S}_1^{\perp T} = (S_1^{\perp T} B_1^\perp)^{-1} S_1^{\perp T}, \quad C_1^T = (B_1^T S_1^*)^{-1} B_1^T. \quad (69)$$

For this choice \tilde{S}_1^\perp is a *maximal* system of linear independent vectors orthogonal to the range of S_1^* . Hence, the size of η , i.e., the number of decoupling

⁷In [13, 14], $B_1 = S_1^*$ and $B_1^\perp = S_1^\perp$ is taken.

PDEs, becomes maximal for this choice. We further get $\mathbf{D}_1 = \mathbf{I}_{N_*}$ (see (18)) and, using (20),

$$\mathbf{R} = \mathbf{A}_{1,het}. \quad (70)$$

Comparing the definitions of the transformed variables for the generalized formulation (21) and the standard formulation in [10, 9, 11] we have the equality

$$\boldsymbol{\xi}_* = (\boldsymbol{\xi}_{sorp}, \boldsymbol{\xi}_{min}, \boldsymbol{\xi}_{kin}). \quad (71)$$

Using all this and the definition of $\mathbf{A}_{1,*}$ (cf. (20)) one sees that the PDEs (22)

$$\partial_t(\theta \boldsymbol{\xi}_*) + L \boldsymbol{\xi}_* = \theta \mathbf{R} \bar{\boldsymbol{\xi}}_{het} + \theta \mathbf{A}_{1,*} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}})$$

are identical to the equations

$$\begin{aligned} \partial_t(\theta \boldsymbol{\xi}_{sorp}) + L \boldsymbol{\xi}_{sorp} &= \theta(\mathbf{r}_{sorp,li} + \mathbf{A}_{1,sorp} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}})) \\ \partial_t(\theta \boldsymbol{\xi}_{min}) + L \boldsymbol{\xi}_{min} &= \theta(\mathbf{r}_{min} + \mathbf{A}_{ld} \mathbf{r}_{sorp,ld} + \mathbf{A}_{1,min} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}})) \\ \partial_t(\theta \boldsymbol{\xi}_{kin}) + L \boldsymbol{\xi}_{kin} &= \theta \mathbf{A}_{1,kin} \mathbf{r}_{kin}(\mathbf{c}, \bar{\mathbf{c}}). \end{aligned}$$

in [9, Eq. (3.19)-(3.21)].

Furthermore, for \mathbf{Y}, \mathbf{Z} defined in (25) the equations

$$\mathbf{Y} = \begin{pmatrix} \mathbf{S}_{1,het}^* & \mathbf{S}_{1,kin}^* \end{pmatrix}, \quad \mathbf{Z} = \mathbf{B}_1^\perp$$

hold. Plugging \mathbf{R} (70) and the partitioning of $\boldsymbol{\xi}_*$ (71) in the definition of the additional variables (39) of the generalized reduction scheme and using the fact that

$$\mathbf{A}_{1,het} = \begin{pmatrix} \mathbf{I}_{J_{sorp,li}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{ld} & \mathbf{I}_{J_{min}} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}$$

in the case of the standard formulation of the reduction scheme (compare (17) with [9, Eq. (3.10)]) gives

$$\begin{pmatrix} \tilde{\boldsymbol{\xi}}_{sorp} \\ \tilde{\boldsymbol{\xi}}_{min} \\ \tilde{\boldsymbol{\xi}}_{kin} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\xi}_{sorp} \\ \boldsymbol{\xi}_{min} \\ \boldsymbol{\xi}_{kin} \end{pmatrix} - \begin{pmatrix} \bar{\boldsymbol{\xi}}_{sorp,li} \\ \mathbf{A}_{ld} \bar{\boldsymbol{\xi}}_{sorp,ld} + \bar{\boldsymbol{\xi}}_{min} \\ \mathbf{0} \end{pmatrix}.$$

Here we find a tiny difference to the standard formulation of the reduction scheme as it is presented in [9, (3.39)] and [10, (18)] where the additional kinetic variables $\tilde{\boldsymbol{\xi}}_{kin}$ do not appear. However, in the general formulation we just doubled these variables $\tilde{\boldsymbol{\xi}}_{kin} = \boldsymbol{\xi}_{kin}$.

6 Numerical Results

In this section we are going to compare the performance of the reformulation in the version which is called the 'standard formulation' in the previous section, to other numerical methods.

For all numerical computations we use the so-called "2-D advective easy test case" of the numerically challenging MoMaS benchmark. See [5] for a description of the benchmark.

In Sec. 6.1 we compare our performance results with those given by other groups. In Sec. 6.2 we try to simulate other numerical methods by our code in order to present results which were obtained using the same numerical kernel for all methods.

6.1 Comparison of the CPU time with other groups

A detailed comparison of the easy test case results of all participants of the benchmark can be found in a synthesis article [4]. So here only a short summary will be given.

For the "2D advective easy test case", three groups presented results. Lagneau and van der Lee presented results computed with the code HYTEC (see [15]), which uses iterative operator splitting (SIA). Mayer and MacQuarrie presented results computed with the code MIN3P (see [16]), which uses the direct substitutional approach (DSA). Hoffmann, Kräutle and Knabner presented results using the presented reformulation method, Sec. 5.2 (for the specific matrices $\mathbf{B}_1, \mathbf{B}_1^\perp, \mathbf{B}_2, \mathbf{B}_2^\perp$ see [10]). In Fig. 1 the normalized CPU times of the three participants in dependence on the number of cells can be seen. This figure is taken from [4].

The other groups give results only for coarser grids (in our opinion too coarse grids, see Sec. 6.3). By extrapolating the lines in Fig. 1 one can estimate that the reduction scheme is faster by a factor greater than five compared with the second fastest code.

The main drawback of DSA is that it leads to a nonlinear system which is difficult to solve numerically because the reaction constants appear in the global system. This is avoided by using the resolution function from Sec. 4. Hence using the resolution function seems to be the preferable method to reduce the number of equations in the global system. Another drawback of DSA is that the decoupling of linear PDEs is not possible. A discussion of the drawbacks of SIA can be found in the following section.

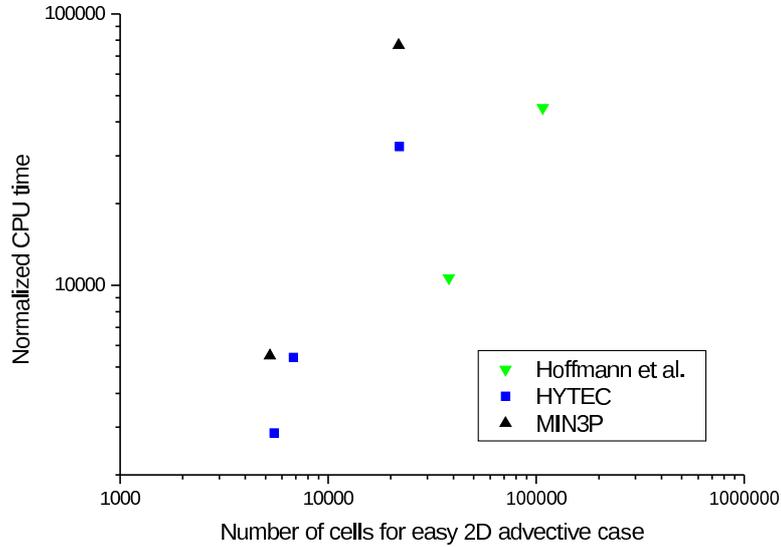


Figure 1: CPU time in the 2D “advective easy test case”

6.2 Comparison of the CPU time with other methods

In this section the universality of the implementation of the reduction scheme (see Sec. 3) is used to compare the CPU times of the reduction scheme with those of the global-ODE approach and those of iterative operator splitting. The time step size is chosen adaptively, based on the number of Newton steps.

In a first comparison the variables and equations of the reduction scheme are used, but the η -equations are added to the global problem and so are solved simultaneously with the global problem, i.e., the decoupling of the η -equations is not exploited by the numerical solver. By doing so one can see the saving of CPU time gained by the decoupling of the η -equations. In Table 2 these results are named “no decoupling of η -equations”.

To obtain the global-ODE approach with the implementation of the generalized reduction scheme one has to choose C_1 according to (67).

To do SIA with the generalized reduction scheme one has to use same matrix C_1 as in global-ODE approach, but has to do some modifications in the implementation. The main thing is to comment out the terms with $D_{\xi_{glob}} \xi_{loc}$ in the Jacobian and to plug equation (45) into (46).

Moreover a stopping criterion for the SIA iteration has to be chosen. For SIA we use the same type of stopping criterion as we use for the reduction

scheme and the global-ODE approach, i.e.,

$$|\mathbf{r}| < \max\{Eps, |\mathbf{r}_0|Red\} \quad (72)$$

where \mathbf{r}_0 is the initial residual and Eps , Red are user specified parameters. However, in order to get a result in finite time, we had to relax the criterion for SIA by multiplying Red by a factor of ten and multiplying Eps by a factor of 100 or 200, i.e., we tried two different values $Eps = 2 \cdot 10^{-8}, 10^{-8}$.

All computations presented are done with the same implementation based on M++ (see [23]), and a FV stabilization with exponential upwinding is used because of the convection dominance.

In Table 2 the CPU time, the number of time steps and the average number of global Newton steps (iteration steps in the case of SIA) are given.

Table 2: CPU time, time steps, Newton steps (for SIA: iteration steps) for different methods

	cells	CPU time	time steps	Newton steps
reduction scheme	26880	6260.5	12921	2.14
reduction scheme	38016	11269.3	14602	2.51
no decoupling of η -equations	26880	10266.0	13123	2.15
no decoupling of η -equations	38016	18498.8	14779	2.52
global-ODE approach	26880	10844.5	13077	2.25
global-ODE approach	38016	19647.9	14986	2.63
SIA ($Eps = 2 \cdot 10^{-8}$)	9504	12031.0	14979	15.3
SIA ($Eps = 10^{-8}$)	9504	14404.3	16726	17.1
SIA ($Eps = 2 \cdot 10^{-8}$)	26880	35437.4	15680	16.6
SIA ($Eps = 10^{-8}$)	26880	42782.3	17949	18.4

Comparing the results of the reduction scheme and the case “no decoupling of η -equations” one can see that the gain of CPU time by the decoupling of the η -equations is a bit more than one third. Furthermore it can be seen, if one does not want to use the reduction scheme, then the global-ODE approach is the best thing one can do.

The SIA approach has two drawbacks. The first one is that the CPU time is much higher than for all the other methods. Compared with the reduction scheme the CPU time is higher by a factor greater than five. The second drawback of the SIA approach is that it is not possible to get

solutions that are as precise as the solutions of the other methods because it is not feasible to choose the same stopping criterion. Using SIA it is necessary to enlarge the stopping parameter Eps , i.e., the error, by a factor of about one hundred, because otherwise the convergence gets so slow that the method is practically unusable. How much more CPU time is needed when the stopping parameter Eps is diminished only by a factor of two is displayed in Table 2.

Analyzing the SIA method one sees that it has a linear convergence behavior like it is expected for a fixed point method. After some iteration steps (about five) in each iteration the residual is reduced by an approximately constant factor. The reason why a severe stopping criterion is practically unusable is that sometimes this factor is only 1.3 depending on the problem and the time step size. Then it would take 18 iteration steps to reduce the residual by a factor of one hundred.

6.3 Grid convergence

The implementation of the reduction scheme also runs on parallel computers. So the use of very fine grids is possible. Furthermore the implementation is done for unstructured grids. So pre-adapted grids can be used. A suitable pre-adapted grid is shown in Fig. 2.

To check the convergence of the method a reference solution on a fine, pre-adapted grid with 608256 cells is computed and the differences between the solutions and the reference solution measured in the discrete L^2 -norm are computed (see Table 3). We observe a numerical convergence order of $h^{2/3}$. Since the given permeability is discontinuous, the water flow has only low regularity. The dispersion tensor depends on this flow. So the solution of the transport problem has very low regularity. Hence, such a low numerical convergence order is not really surprising.

Table 3: L^2 -error for different grids with respect to the reference solution

number of cells	pre-adapted	error in discrete L^2 -norm
26880	no	0.0120
107520	no	0.0076
38016	yes	0.0100

Most streamlines go through the high velocity zone ($1 \leq x \leq 1.1$, $0.9 \leq$

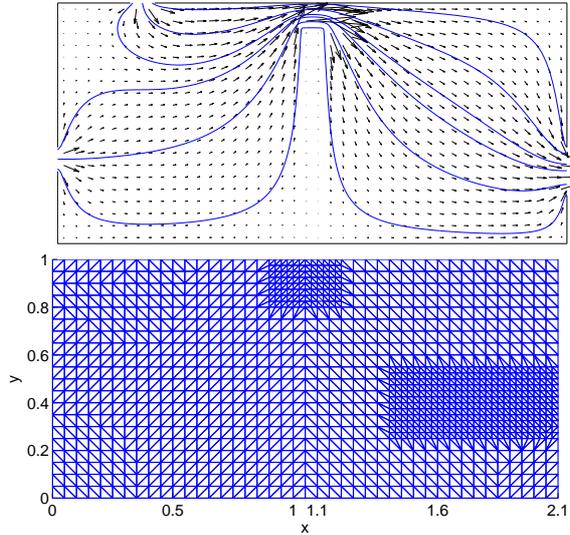


Figure 2: Water flow (top) and grid refined in the high velocity zone and near the outflow (bottom)

$y \leq 1.0$) (see Fig. 2). Hence, the number of cells in this part of the grid is crucial for the accuracy of the solution in the right half of the domain. The numbers of cells in y -direction in the high velocity zone are 64 for the reference solution, 8 for the grid with 26880, 7 for the grid used by HYTEC, 10 for the grid used by MIN3P (in the case 'without CPU time limitation'; in the case 'with CPU time limitation' a coarser grid was used; see [4] for the grids used by HYTEC and MIN3P), 4 for the grid used by GDAE2D (see [7]). The grid used to compute the reference solution is the only pre-adapted grid in this list. All other grids are regular grids.

The grid with 26880 cells is comparable with the grids used by HYTEC and MIN3P. There are considerable differences between the solution on the grid with 26880 cells and the reference solution, especially between $x=1$ and 1.2 at the bottom, and concerning the extent of the red zone on top right (see Fig. 3). Hence, a mesh with this seems too coarse to give good results on the right-hand side of the domain.

The possibility to use a sufficiently fine mesh (our reference solution) we owe to the high computational efficiency of the reduction method.

So in the "2D advective test case" of the MoMaS benchmark the implementation of the reduction scheme was the only code able to compute a

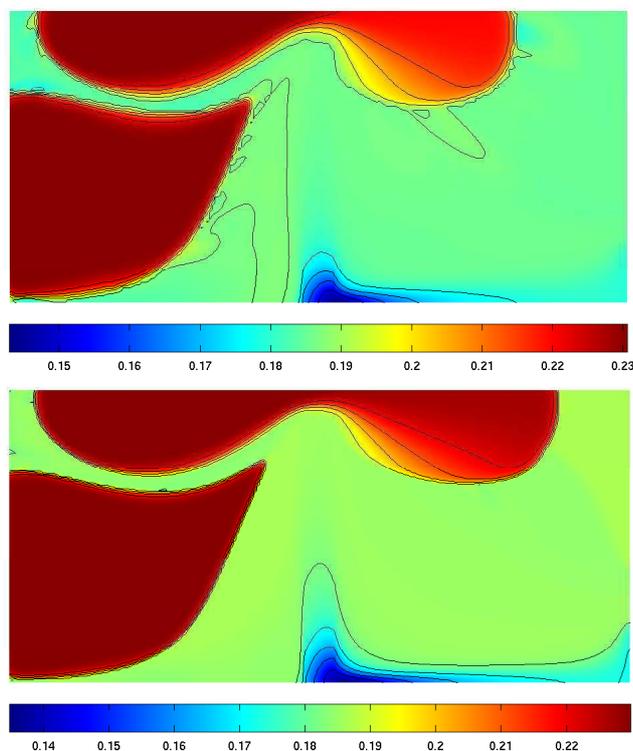


Figure 3: Component X_3 at time 1000. Grid with 26880 cells (top), reference solution with 608256 cells (bottom).

solution on a sufficiently fine mesh, up to now.

7 Conclusions

The reduction scheme decreases the number of coupled equations and with it the CPU time significantly. It is possible to apply the method to realistic problems with large reaction constants and concentrations ranging over many orders of magnitude. The numerically challenging MoMaS benchmark was solved successfully. In this benchmark the reduction scheme is more than five times faster than the software of the second fastest group. Compared with the standard methods SIA and DSA, the reduction scheme is more efficient for this type of problem. Due to the efficiency of the reduction scheme, simulations on finer grids are accomplishable even on standard PCs.

For realistic 2D problems such fine grids are necessary to get an accurate solution.

References

- [1] Amir, L., Kern, M.: A global method for coupling transport with chemistry in heterogeneous porous media, *Comp. Geosci.* 14, 465481 (2010), doi:10.1007/s10596-009-9162-x
- [2] Bethke, C.M.: *Geochemical reaction modeling*, Oxford University Press, New York (1996)
- [3] Carrayrou, J.: *Modélisation du transport de solutés réactifs en milieu poreux saturé*, doctoral thesis, Université Louis Pasteur Straßburg (2001)
- [4] Carrayrou, J., Hoffmann, J., Knabner, P., Kräutle, S., de Dieuleveult, C., Erhel, J., Van der Lee, J., Lagneau, V., Mayer, K.U., McQuarrie, K.T.B.: Comparison of numerical methods for simulating strongly non-linear and heterogeneous reactive transport problems — the MoMaS benchmark case, *Comp. Geosci.* 14, 483-502 (2010), doi:10.1007/s10596-010-9178-2
- [5] Carrayrou, J., Kern, M., Knabner, P.: Reactive transport benchmark of MoMaS, *Comp. Geosci.* 14, 385-392 (2010), doi:10.1007/s10596-009-9157-7
- [6] Carrayrou, J., Mosé, R., Behra, P.: Operator-splitting procedures for reactive transport and comparison of mass balance errors, *J. Cont. Hydr.* 68, 239-268 (2004)
- [7] de Dieuleveult, C.: *Un modèle numérique global et performant pour le couplage géochimie-transport*, doctoral thesis, Université de Rennes 1 (2008)
- [8] de Dieuleveult, C., Erhel, J., Kern, M.: A global strategy for solving reactive transport equations, *J. Comput. Phys.* 228, 6395-6410 (2009)
- [9] Hoffmann, J.: *Reactive transport and mineral dissolution/precipitation in porous media: efficient solution algorithms, benchmark computations and existence of global solutions*, doctoral thesis, University Erlangen-Nuremberg (2010)

- [10] Hoffmann, J., Krättele, S., Knabner, P.: A parallel global-implicit 2-d solver for reactive transport problems in porous media based on a reduction scheme and its application to the MoMaS benchmark problem, *Comp. Geosci.* 14, 421-433 (2010), doi:10.1007/s10596-009-9173-7
- [11] Krättele, S.: General multi-species reactive transport problems in porous media: efficient numerical approaches and existence of global solutions, habilitation thesis, University Erlangen–Nuremberg (2008)
- [12] Krättele, S.: The semismooth Newton method for multicomponent reactive transport with minerals, *Advances Water Res.* 34, 137-151 (2011), doi:10.1016/j.advwatres.2010.10.004
- [13] Krättele, S., Knabner, P.: A new numerical reduction scheme for fully coupled multicomponent transport-reaction problems in porous media, *Water Resour. Res.* 41 (2005), W09414, doi:10.1029/2004WR003624
- [14] Krättele, S., Knabner, P.: A reduction scheme for coupled multicomponent transport-reaction problems in porous media: Generalization to problems with heterogeneous equilibrium reactions, *Water Resour. Res.* 43 (2007), W03429, doi:10.1029/2005WR004465
- [15] Lagneau, V., van der Lee, J.: HYTEC results of the MoMaS reactive transport benchmark, *Comp. Geosci.* 14, 435-449 (2010), doi:10.1007/s10596-009-9159-5
- [16] Mayer, K.U., MacQuarrie, K.T.B.: Formulation of the multicomponent reactive transport code MIN3P and implementation of MoMaS benchmark problems, *Comp. Geosci.* 14, 405-419 (2010), doi:10.1007/s10596-009-9158-6
- [17] Mayer, K.U., Frind, E.O., Blowes, D.W.: Multicomponent reactive transport modeling in variably saturated porous media using a generalized formulation for kinetically controlled reactions, *Water Resour. Res.* 38, 1174-1194 (2002), doi:10.1029/2001WR000862
- [18] Morel, F., Hering, J.: Principles and applications of aquatic chemistry, Wiley, New York (1993)
- [19] Saaltink, M.W., Ayora, C., Carrera, J.: A mathematical formulation for reactive transport that eliminates mineral concentrations, *Water Resour. Res.* 34, 1649-1656 (1998)

- [20] Saaltink, M.W., Carrera, J., Ayora, C.: A comparison of two approaches for reactive transport modelling, *J. Geochem. Explo.* 69-70, 97-101 (2000)
- [21] Scheidegger, A.E.: General theory of dispersion in porous media, *J. Geophys. Res.* 66, 3273-3278 (1961)
- [22] van der Lee, J., De Windt, L., Lagneau, V., Goblet, P.: Modul-oriented modeling of reactive transport with HYTEC, *Computers & Geosciences* 29, 265-275 (2003)
- [23] Wieners, C.: Distributed point objects. A new concept for parallel finite elements, in *Domain decomposition methods in science and engineering, Lecture notes in computational science and engineering*, vol. 40, R. Kornhuber, R. Hoppe, J. Piaux, O. Pironneau, O. Widlund, J. Xu (editors), Springer, 175-183 (2004)

PREPRINTS

DES INSTITUTS FÜR ANGEWANDTE MATHEMATIK DER UNIVERSITÄT ERLANGEN-NÜRNBERG

ZULETZT ERSCHIENENE BEITRÄGE:

- 330 J. HASLINGER , G. LEUGERING , M. KOČVARA , M. STINGL: *Multidisciplinary Free Material Optimization*
- 331 B. SCHMIDT , M. STINGL , D. A. BERRY , M. DÖLLINGER: *Material parameter optimization in a multi-layered vocal fold model*
- 332 M. KAISER, A. THEKALE: *Solving nonlinear feasibility problems with expensive functions*
- 333 I. BOMZE , G. EICHPELDER: *Copositivity detection by difference-of-convex decomposition and ω -subdivision*
- 334 M. PRECHTEL, G. LEUGERING, P. STEINMANN, M. STINGL: *Towards optimization of crack resistance of composite materials by adjustment of fiber shapes*
- 335 A. M. KHLUDNEV , G. LEUGERING: *Optimal control of cracks in elastic bodies with thin rigid inclusions*
- 336 M. PRECHTEL, P. LEIVA RONDA, R. JANISCH, A. HARTMAIER, G. LEUGERING, P. STEINMANN, M. STINGL: *Simulation of fracture in heterogeneous elastic materials with cohesive zone models*
- 337 E. MARCHAND: *Combined Deterministic-Stochastic Sensitivity Analysis; Application to Uncertainty Analysis.*
- 338 G. EICHPELDER , T.X.D. HA: *Optimality conditions for vector optimization problems with variable ordering structures*
- 339 F.A. RADU, N. SUCIU, J. HOFFMANN, A. VOGEL, O. KOLDITZ, C-H. PARK , S. ATTINGER: *Accuracy of numerical simulations of contaminant transport in heterogeneous aquifers: a comparative study*
- 340 M. KAISER, K. KLAMROTH, A. THEKALE: *Test examples for nonlinear feasibility problems with expensive functions*
- 341 J. JAHN, T.X.D. HA: *New Order Relations in Set Optimization*
- 342 G. EICHPELDER , J. POVH: *On reformulations of nonconvex quadratic programs over convex cones by set-semidefinite constraints*
- 343 T.X.D. HA AND J. JAHN: *Properties of Bishop-Phelps Cones*
- 344 N. RAY, CH. ECK, A. MUNTEAN, P. KNABNER: *Variable Choices of Scaling in the Homogenization of a Nernst-Planck-Poisson Problem*
- 345 A. MUNTEAN , T. L. VAN NOORDEN: *Corrector estimates for the homogenization of a locally-periodic medium with areas of low and high diffusivity*
- 346 N.SUCIU , S.ATTINGER , F.A.RADU , C.VAMOŞ , J.VANDERBORGH , H.VERECKEN , P. KNABNER: *Solute transport in aquifers with evolving scale heterogeneity*
- 347 F. BRUNNER, F. A. RADU, M. BAUSE, P. KNABNER : *Optimal order convergence of a modified BDM_1 mixed finite element scheme for reactive transport in porous media*
- 348 G. EICHPELDER: *Cone-valued maps in optimization*
- 349 G. EICHPELDER , J. POVH: *On the set-semidefinite representation of nonconvex quadratic programs over arbitrary feasible sets*
- 350 S. KRÄUTLE: *Existence of global solutions of multicomponent reactive transport problems with mass action kinetics in porous media*
- 351 N. RAY, T. VAN NOORDEN, F. FRANK, P. KNABNER: *Colloid and Fluid Dynamics in Porous Media including an Evolving Microstructure*
- 352 F. FRANK, N. RAY, P. KNABNER: *Numerical Investigation of a Homogenized Stokes-Nernst-Planck-Poisson Problem*