



DecTree v1.0 - Chemistry speedup in reactive transport simulations: purely data-driven and physics-based surrogates

Marco De Lucia¹ and Michael Kühn^{1,2}

¹GFZ German Research Centre for Geosciences, Telegrafenberg, 14473 Potsdam, Germany

²University of Potsdam, Institute of Geosciences, Hydrogeology, Potsdam, Germany

Correspondence: M. De Lucia (delucia@gfz-potsdam.de).

Abstract. The computational costs associated with coupled reactive transport simulations are mostly due to the chemical subsystem: replacing it with a pre-trained statistical surrogate is a promising strategy to achieve decisive speedups at the price of small accuracy losses and thus to extend the scale of problems which can be handled. We introduce a hierarchical coupling scheme in which “full physics”, equation-based geochemical simulations are partially replaced by surrogates. Errors on mass balance resulting from multivariate surrogate predictions effectively assess the accuracy of multivariate regressions at runtime: inaccurate surrogate predictions are rejected and the more expensive equation-based simulations are run instead. Gradient boosting regressors such as xgboost, not requiring data standardization and being able to handle Tweedie distributions, proved to be a suitable emulator. Finally, we devise a surrogate approach based on geochemical knowledge, which overcomes the issue of robustness when encountering previously unseen data, and which can serve as basis for further development of hybrid physics-AI modelling.

1 Introduction

Coupled reactive transport simulations (Steeffel et al., 2005, 2015) are very expensive, effectively hampering their wide applications. While hydrodynamic simulations on finely resolved spatial discretisations, containing millions of grid elements, are routinely run on common workstations, the order of magnitude of the computationally affordable reactive transport simulations on the same hardware decreases by a factor of ten to a hundred as soon as chemical reactions are coupled in (De Lucia et al., 2015; Jatnieks et al., 2016; De Lucia et al., 2017; Leal et al., 2020; Prasianakis et al., 2020). This usually requires oversimplifications of the subsurface domain, reduced to 2D or very coarse 3D, and of the geochemical complexity as well.

In the classical *operator splitting* such as Sequential Non-Iterative Approach (SNIA), the three interacting physical processes hydrodynamic flow, solutes transport and chemical interactions between solute species and rock forming minerals are solved sequentially. Chemistry usually represents the bottleneck for coupled simulations with up to 90% of computational time (De Lucia and Kühn, 2013; De Lucia et al., 2015; Leal et al., 2020). The numerical model for geochemical speciation and reactions requires in general the integration of one stiff differential-algebraic system of equations per grid element per



simulation time step. Parallelisation is thus required to tackle large spatial discretisations, which is why many modern codes
25 are developed to run on high performance computing (HPC) clusters with many thousand of CPUs. However, parallelisation
alone still is not sufficient to ensure numerical convergence of the simulations, a problem routinely encountered by many
practitioners. Furthermore, large uncertainties affect the phenomenological model itself (i.e., reaction kinetics spanning over
orders of magnitude; consistent activity model for concentrated solutions usually encountered in the subsurface, . . .) and even
larger concern the parametrisation of the subsurface, regarding the spatial heterogeneity of rocks, which is mostly unknown
30 and hence often disregarded (De Lucia et al., 2011). It thus may appear unjustified to allocate large computational resources to
solve very expensive, yet still actually oversimplified or uncertain problems. Removing the computational cost associated with
reactive transport modelling is thus of paramount importance to ensure its wide application on a range of otherwise practically
unfeasible problems (Prommer et al., 2019).

The much desired speedup of this class of numerical models has been the focus of intensive research in the last few years.
35 Among the proposed solutions, Jatnieks et al. (2016) suggests to replace the “full physics” numerical models of the geochem-
ical subsystem with emulators or surrogates, employed at runtime during the coupled simulations. A surrogate in this sense is
a statistical multivariate regressors which has to be trained in advance on a set of pre-calculated “full physics” solutions of the
geochemical model at hand, spanning the whole parameter range expected for the simulations. Since the regressors are much
quicker to compute than the setup and integration of a differential algebraic system of equations (DAE), this promises a signif-
40 icant speedup and has thus found resonance in the scientific community (e.g., De Lucia et al., 2017; Laloy and Jacques, 2019;
Guérrillot and Bruyelle, 2020). However, all approximations and especially purely data-driven surrogates introduce accuracy
losses into the coupled simulations. These must be kept low in order to generate meaningful simulation results. Ultimately,
replacing a fully fledged geochemical simulator with surrogate equals to trading computational time for accuracy of the sim-
ulations. Due to the non-linear nature of geochemical sub-processes, even small errors in surrogate predictions propagate in
45 successive iterations so that diverging trajectories for the coupled models originate from only few time-steps, leading to un-
physical results. Mass and charge imbalances, i.e., “creation” of matter, happen to be the most common source of unphysicality
in our early tests. It is thus of paramount importance to obtain highly accurate surrogates, which in turn may require very large
and densely sampled training datasets and training times.

Any regression algorithm can be employed to replace the “full physics” equation-based geochemical model. The thriving
50 developments in data science and machine learning in the recent years have produced many different and efficiently imple-
mented regressors readily available and usable in high-level programming languages such as python or R. Among the most
known ones there are gaussian processes, support vector machines, artificial neural networks and decision-tree based algo-
rithms such as random forest or gradient boosting. Most of these algorithms are “black boxes”, which non-linearly relate many
output variables to many input variables. Their overall accuracy can be statistically assessed by measuring their performances
55 on the training dataset or on a subset of the available training data left out for the specific purpose of testing the models. In any
case these training and/or test datasets must be obtained beforehand computing an appropriate number of points with the “full
physics” model. Geochemistry is usually largely multivariate, meaning that many input and many output variables are passed
to and from the geochemical subsystem at each time step. In general, different regressors may capture each output variable in



better fashion depending on many factors (e.g., the problem at hand, where variables display different non-linear behaviors; the sampling density of training dataset, which may be biased). For this reason, we argue that the most sensible choice for a surrogate modelling framework is that of *multiple multivariate regression*: one multivariate regressor - making use of many or all inputs as predictors - is trained independently for each distinct output variable, while the choice of regressor may vary from variable to variable. With algorithms such as Artificial Neural Networks (ANN) it is possible to train one single network, and hence one single surrogate, for all variables at once. However, ANN usually require long CPU times for training and quite large training datasets.

In praxis, the CPU-time, the user interactions and the overall skills required for optimally training complex regressors cannot be underestimated and may prove overwhelming for a geoscientist. The whole process of hyperparameter tuning, required by most advanced machine learning algorithms, while being an active area of development and research, is still hardly fully automatable.

This work showcases and analyses two different approaches for surrogate geochemical modelling in reactive transport simulations. The first is completely data-driven, disregarding any possible knowledge about the ongoing process. In the second approach, we derive a surrogate which exploits the actual equations solved by the full physics representation of chemistry. Both are applied and evaluated on the same 1D benchmark implemented in a simple reactive transport framework. Our implementation of coupled reactive transport includes a hierarchical submodel coupling strategy, which is advantageous when different accuracy levels for the predictions of one sub-process are available.

2 Methods: simulation environment and benchmark problem

The versioned R code used for DecTree v.1.0 model setup and evaluation is referenced in section “code availability”. It is based on version v0.0.4 of the `RedModRphree` package for the R environment (R Core Team, 2020), which is also referenced in section “code availability”. It makes use of the geochemical simulator PHREEQC (Appelo et al., 2013). `RedModRphree` supersedes the in-house developed R-PHREEQC interface `Rphree` (<https://rphree.r-forge.r-project.org/>, De Lucia and Kühn, 2013).

The benchmarks and the performance measurements refer to computations run on a recent desktop workstation equipped with Intel Xeon W-2133 CPU with clock at 3.60 GHz and DDR4 RAM at 2.666 GHz under Linux kernel 5.9.14 and R version 4.0.3. If not otherwise specified, only one CPU core is employed for all computational tasks. Since chemistry is inherently an embarrassing parallel task, the speedup achieved on a single CPU as in this work will transfer - given large enough simulation grids making up for the overhead - on parallel computations.

2.1 Numerical simulation of flow and transport

We consider a stationary, fully-saturated, incompressible, isothermal 1D Darcy flow in a homogeneous medium. Transport is restricted to pure advection and the feedback of mineral precipitation and dissolution on porosity and permeability is also disregarded; the fluid density is also considered constant. Advection is numerically computed via a *forward Euler* explicit



resolution scheme:

$$C_i(x, t + 1) = C_i(x, t) - u \cdot \Delta t \frac{C_i(x, t) - C_i(x - 1, t)}{\Delta x} \quad (1)$$

where u is the module of Darcy velocity, $C_i(x, t)$ the volumetric concentration (molality) of the i -th solute species at the point x and time t , and Δx the size of a grid element. For this scheme, the Courant-Friedrichs-Lewy stability condition (CFL)

95 imposes that the Courant number ν be less than or equal to 1:

$$\nu = \frac{u \cdot \Delta t}{\Delta x} \leq 1 \quad (2)$$

For Courant numbers less than 1, numerical dispersion arises; the scheme is unstable for $\nu > 1$. The only both stable and precise solution for advection is with $\nu = 1$. Thus, the CFL condition is very limiting in Δt : a factor two refinement in the spatial discretisation corresponds to a factor two decrease in Δt , thus obtaining double required iterations. Note that equation 1 is not
100 written in terms of porosity, so that effectively the Darcy velocity is assumed equal to the fluid flux or, alternatively, porosity is equal to unity. This assumption does not have any impact on the calculations besides the initial scaling of the system.

The implemented advection relies on transport of total elemental concentrations instead of the actual dissolved species, since all solutes are subjected to the very same advection equation (Parkhurst and Wissmeier, 2015). Special attention needs to be paid to pH which is defined in terms of activity of protons:

105 $\text{pH} = -\log_{10}([\text{H}^+])$

and is hence not additive. Assuming that the activity coefficient for protons is constant throughout the simulation, the activity $[\text{H}^+]$ can be actually transported instead of pH. Charge imbalance and redox potential (pe) can be safely disregarded for this redox-insensitive model, with absolutely insignificant errors when compared to the same problem simulated, e.g., with PHREEQC's ADVECTION keyword (not shown).

110 2.2 The chemical benchmark

The chemical benchmark used throughout this work is inspired by Engesgaard and Kipp (1992) and is well known, with many variants, in the reactive transport community (e.g., Shao et al., 2009; Leal et al., 2020). It was chosen since it has been studied by many different authors and it is challenging enough from a computational point of view.

At the inlet of a column, conventionally on the left side in the pictures throughout this work, a 0.001 molal magnesium chloride (MgCl_2) solution is injected into a porous medium whose initial solution is at thermodynamic equilibrium with calcite.
115 With the movement of the reactive front, calcite starts to dissolve and dolomite is transiently precipitated. Kinetic control is imposed on all mineral reactions following a Lasaga rate expression from Palandri and Kharaka (2004), limited to only neutral and H^+ mechanisms (parameters are summarised in Table 1) and constant reactive surfaces, hence independent on the actual amounts of minerals. Precipitation rate - relevant only for dolomite - is set equal to the rate of dissolution. Temperature is set
120 for simplicity at constant 25°C in disregard to actual physical meaningfulness of the model concerning dolomite precipitation (Möller and De Lucia, 2020). Detailed initial and boundary conditions are summarized in Table 2. To achieve a complete



Table 1. Parameters for kinetic control for dissolution and precipitation of calcite and dolomite. k is given in $\text{mol m}^{-2} \text{s}^{-1}$, E_a in kJ mol^{-1} and reactive surface in $\text{m}^2/\text{kgH}_2\text{O}$.

Mineral	H^+ mechanism			Neutral mechanism		
	log k	E_a	H^+ order	log k	E_a	reactive surface
calcite	-0.30	14.4	1	-5.81	23.5	3.20
dolomite	-3.19	36.1	0.5	-7.53	52.2	0.32

Table 2. Initial (IC) and boundary (BC) conditions for the benchmark problem.

	C	Ca	Cl	Mg	pH	calcite	dolomite
	molal	molal	molal	molal	-	mol	mol
IC	$1.2279 \cdot 10^{-4}$	$1.2279 \cdot 10^{-4}$	0.00	0.00	9.91	$2.07 \cdot 10^{-3}$	0.00
BC	0.00	0.00	$0.2 \cdot 10^{-2}$	$0.1 \cdot 10^{-2}$	7		

description of the chemical system at any time, seven input variables are required: pH, C, Ca, Mg, Cl, calcite and dolomite - those can be considered, with an abuse of language, *state variables*, in the sense that they constitute the necessary and sufficient inputs of the geochemical subsystem, and all reactions only depend on them. The outcome of the “full physics” calculations is completely defined (at least with the simplifications discussed above) by four distinct quantities: the amounts of reaction affecting the two minerals calcite (i) and dolomite (ii) in the given time step, from which the changes in solutes Ca, Mg and C can be backcalculated; Cl (iii), which is actually non-reactive; and pH (iv). In a completely process-agnostic, data-driven framework, however, the relationships between minerals and aqueous concentrations are disregarded, and the output of the chemical subsystem is expressed solely in terms of the input variables.

130 2.3 Reference simulations and training data

For the remainder of this work, the geochemical benchmark described above is solved on a 1D column of length 0.5 m, with constant fluid velocity of $u = 9.375 \cdot 10^{-6}$ m/s. The domain is discretised with grid refinements ranging between 50 and 500 grid elements. Higher refinements have a double effect: on one side larger grids obviously increase the overall computational load, and in particular for chemistry; on the other side, given the restriction of the implemented forward Euler explicit advection scheme, the time stepping required for the coupled simulations in order to be free of numerical dispersion decreases accordingly. Smaller time steps decrease the computational load for geochemistry for each iteration, since they require shorter time integrations, but also require more coupled iterations to reach the same simulation time. More iterations mean also that there are more chances for errors introduced by surrogates to further propagate into the simulations in both space and time. In presence of significant overhead due to, e.g., data passing between different softwares or the setup of geochemical simulations,



140 the advantage due to shorter time steps vanishes. However, these aspects become more relevant in the context of parallelisation of geochemistry and are not addressed in the present work.

All coupled simulations, both reference (full physics) and with surrogates, are run with constant time step either honouring the CFL condition with $\nu = 1$, and thus free of numerical dispersion, or, when assessing how the speedup scales with larger grids, a fixed time step small enough for the CFL condition (eq. 2) to be satisfied for every discretisation. As previously noted, 145 the resulting simulations will be affected by grid-dependent numerical dispersion, which we do not account for in the present work. This makes the results incomparable in terms of transport across grids. However, since the focus is on the acceleration of geochemistry through pre-computed surrogates, this is an acceptable simplification.

The comparison between the reference simulations, obtained by coupling of transport with the PHREEQC simulator, and those obtained with surrogates is based on an error measure composed as the geometric mean of the relative RMSEs of each 150 variable i , using the variable's maximum at a given time step as norm:

$$\text{Error}_t = \exp \left\{ \frac{1}{m} \sum_i^m \ln \frac{\sqrt{\frac{1}{n} \sum_j^n (\text{ref}_{i,j} - \text{pred}_{i,j})^2}}{\max_t(\text{pred}_i)} \right\} \quad (3)$$

where m is the number of variables to compare, n the grid dimension and t the particular time step where the error is computed.

In this work the datasets used for training the surrogates are obtained directly by storing all calls to the full physics simulator and its responses in the reference coupled reactive transport simulations, possibly limited to a given simulation time. This way 155 of proceeding is considered more practical than, e.g., an a priori sampling of a given parameter space, where the bounds of the parameter space are defined by the ranges of the input/output variables occurring in the reference coupled simulations. This strategy mimics the problem of wanting to train a surrogates directly at runtime during the coupled simulations. Furthermore, an a priori, statistical sampling of parameter space, in absence of restrictions based on the physical relationships between the variables, would include unphysical and irrelevant input combinations. By employing only the input/outputs tables actually 160 required by the full physics simulations, this issue is automatically solved; however, the resulting datasets will be in general skewed, multimodal and highly inhomogeneously distributed within the parameter space, with highly dense samples in some regions and even larger empty ones.

2.4 Hierarchical coupling of chemistry

In this work we consider only a Sequential Non-Iterative Approach (SNIA) coupling scheme, meaning that the sub-processes 165 flow, transport and chemistry are solved numerically one after another before advancing to the next simulation step. For sake of simplicity, we let the CFL condition (2) for advection dictate the allowable time step for the coupled simulations.

Replacing the time-consuming, equation-based numerical simulator for geochemistry with an approximated but quick surrogate introduces inaccuracies into the coupled simulations. These may quickly propagate in space and time during the coupled simulations and lead to ultimately incongruent and unusable results.



170 A way to mitigate error propagation, and thus to reduce the accuracy required of the surrogates, is represented by a *hierarchy of models* used to compute chemistry at each time step during the coupled simulations. The idea is to first ask the surrogate for predictions, then identify unphysical or unphysical ones, and finally run the full physics chemical simulator for the rejected ones. This way, the surrogate can be tuned to capture with good accuracy the bulk of the training data, and no particular attention needs to be paid to the most difficult “corner cases”. For the highly non-linear systems usually encountered in
175 geochemistry, this is of great advantage. In practice, however, there still is need to have a reliable and cheap error estimation of surrogate predictions at runtime.

It is important to understand that the criteria employed to accept or reject the surrogate predictions depend strictly on the architecture of the multivariate surrogate and on the actual regression method used. Methods such as kriging offer error models, based, e.g., on the distance of the estimated point from the nearest training data. However, in the general case, any
180 error estimation requires first the training and then the evaluation at runtime of a second “hidden” model. Both steps can be time-consuming; furthermore, in the general case one can only guarantee that the error is *expected* - in a probabilistic sense - to be lower than a given threshold.

In a completely data-driven surrogate approach, where each of the output variables is independently approximated by a different multivariate regressor, checking mass conservation is a very inexpensive way to estimate the reliability of a given
185 surrogate prediction, since it only requires the evaluation of linear combinations across predictors and predictions. Other constraints may be added, suited to the chemical problem at hand, such as charge balance. However we only use mass balance in the present work. Figure 1 illustrates this simple hierarchical coupling schematically. For the chemical benchmark of sec-

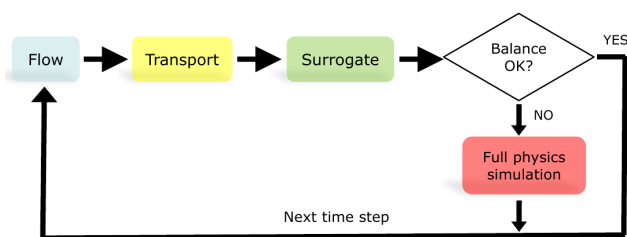


Figure 1. Schematic view of hierarchical sequential non iterative coupling. The decision whether to accept or not the predictions of a multiple multivariate surrogate is based on computing the mass balances for the three elements forming dolomite and calcite before and after reaction, and computing their mean absolute error. If this error exceeds a given threshold, the more expensive equation based geochemical simulator is run instead.

tion 2.2, three mass balance equations can be written, one for each element C, Ca and Mg, accounting for the stoichiometry of the minerals’ brute formulas. If a surrogate prediction exceeds a given, predetermined, tolerance on the Mean Absolute Error
190 of the balance equations, that particular prediction is rejected and a more expensive full physics simulation is run instead.

This approach moderates the need for extremely accurate regressions, especially in instances of non-linear behaviour of the chemical models, for example when a mineral precipitates for the first time or when it is completely depleted, which are hard things for regressors to capture. However, the number of rejected simulations must be low to produce relevant speedups; it is



effectively a trade-off between the accuracy of the surrogates (and efforts and time which goes into it) and the speedup achieved
195 in coupled simulations.

3 Fully data-driven approach

The first approach is a completely general one, fully data-driven and thus process-agnostic: it can be employed for any kind of
numerical model or process which can be expressed in form of input and output tables. In our case, the tables produced by the
geochemical sub-process during the reference coupled simulations are used to train seven multiple multivariate regressors, one
200 for each output.

The reference simulations, and hence the dataset for training the surrogate, are fully coupled simulations on grid 50, 100
and 200 with a fixed time step of 210 s, run until 33600 s or else 161 total coupling iterations. As previously noted, these
simulations are then not comparable among themselves due to significant numerical dispersion; however, from the point of
view of geochemical processes, this strategy has the advantage of spreading the “perturbations” due to transport of the results
205 of geochemistry in the previous time step. Instead of the usual random split of the dataset in train and test subsets, customary
in the machine learning community, we retained only the data resulting from the first 101 iterations for training the surrogates,
and evaluated the resulting reactive transport simulations until iteration 161, where the geochemical surrogate is faced with 60
iterations on unseen or “out of sample” geochemical data. The training dataset comprises tables with 13959 unique rows or
input combinations. All simulations, the reference and with surrogates, are run on a single CPU core.

210 The choice of the regressor for each output is actually arbitrary, and nothing forbids to have different regressors for each
variables, or even different regressors in different regions of parameter space of each variable. Without going into details on all
kinds of algorithms that we tested, we found that decision-tree based methods such as Random Forest and their recent gradient
boosting evolutions appear the most flexible and successful for our purposes. Their edge can in our opinion be resumed by:
(1) implicit feature selection by construction, meaning that the algorithm automatically recognizes which input variables are
215 most important for the estimation of the output; (2) no need for standardisation of both inputs and outputs; (3) ability to
deal with any probability distributions; (4) fairly quick to train with sensible hyperparameter defaults; (5) extremely efficient
implementations available.

The points number (2)-(4) cannot be overlooked. Data normalisation or standardisation techniques, also called “preprocess-
ing” in machine learning lingo, are redundant with decision-tree based algorithms, whereas they have a significant impact on
220 results and training efficiency with other regressors such as Support Vector Machines and Artificial Neural Networks. The
distributions displayed by the variables in the geochemical data are extremely variable and cannot be assumed uniform, gaus-
sian or lognormal in general. We found out that the Tweedie distribution is suited to reproduce many of the variables in the
training dataset. The Tweedie distribution is a special case of exponential dispersion models introduced by Tweedie (1984) and
thoroughly described by Jørgensen (1987), which finds application in many actuarial and signal processing processes (Hassine
225 et al., 2017). A Random Variable Y is a Tweedie distribution of parameter p if $Y \geq 0$, $E[Y] = \mu$ and $Var(Y) = \sigma^2 \mu^p$. This
means that it is a *family* depending on p : gaussian if $p = 0$, Poisson if $p = 1$, gamma if $p = 2$ and inverse gaussian if $p = 3$. The



interesting case, which is normally referred to when using the term “Tweedie”, is when $1 \leq p \leq 2$. This distribution represents positive variables with *positive mass at zero*: meaning that this distribution preserves the “physical meaning” of zero. It’s intuitively an important property when modelling solute concentrations and mineral abundances: the geochemical system solved by the full physics simulator is radically different when, e.g., a mineral is present or not.

Extreme gradient boosting `xgboost` (Chen and Guestrin, 2016) is a decision-tree based algorithm which enjoys an enormous success in the machine learning community in recent years. It has out-of-the-box the capability to perform regression of Tweedie variables and it is extremely efficient in both training and prediction. Using the target Tweedie regression with fixed $p = 1.2$, max tree depth of 20, the default $\eta = 0.3$ and 1000 boosting iterations with early stopping at 50, all results in the dataset are reproduced with great accuracy and the training itself takes around 20 seconds for all seven outputs on our workstation, using four cores. Contrary to the expectation, the accuracy of the predictions is largely enhanced if the labels (i.e., the output table) are scaled before training. We used the max value of each label divided by $1 \cdot 10^{-5}$ as scale. The default evaluation metric when performing Tweedie regression is the Root Mean Squared Log Error:

$$\text{rmsle} = \sqrt{\frac{1}{N} [\ln(\text{pred} + 1) - \ln(\text{label} + 1)]^2} \quad (4)$$

In the previous section it was claimed that in the framework of hierarchical coupling there is no practical need to further refine the regressions. This could be achieved by hyperparameter tuning and by using a different and more adapted probability distribution for each label including proper fitting of parameter p for the Tweedie variables. While this would be of course beneficial, we proceed now by plugging such a rough surrogate into the reactive transport simulations. The coupled simulations with surrogates are performed on the three grids for 161 iterations, setting the tolerance on mass balance to 10^{-5} , 10^{-6} and only relying on the surrogate, meaning with no call to PHREEQC even if a large mass balance error is detected.

In Figure 2 are exemplarily displayed the variables profiles for grid 100 and tolerance 10^{-6} at two different time steps, iteration 101, which is the last one within the training dataset, and at the end of the simulation time, after 60 coupling iterations in “unseen territory” for the surrogates. The accuracy of the surrogate simulations is excellent for the 101st iteration, but by iteration 161, while still acceptable, some discrepancies start to show. The number of rejected surrogate responses at each time step does not remain constant during the simulations, but increases steadily. An overview of all the simulations is given in Figure 3 (top frame). The more stringent mass balance tolerance of 10^{-6} (solid lines) rejects obviously many more simulations which goes hand in hand with the excellent accuracy of the results (Figure 3, bottom panel; error measured with formula of equation 3 excluding pH). It was expected, and it is demonstrated by the evaluation, that starting with the first “out of sample” time step the accuracy of the surrogates significantly drops, which triggers a steep increase of rejected predictions and conversely of calls to PHREEQC. The hierarchical coupling ensures that the errors in the surrogate simulations do not follow the same steep increase, but from this moment on there is a loss of computational efficiency, visible in the simulations with tolerance 10^{-6} , which makes the whole surrogate predictions actually useless in terms of speedup even before making them so inaccurate to be useless. It is also apparent from the error panel in Figure 3 (bottom) that errors introduced in the coupled simulations at early time steps propagate through the rest of the simulations, so that the overall discrepancy between reference and surrogate simulations also steadily increases. Note that this “diverging behavior” also tends to bring the geochemistry

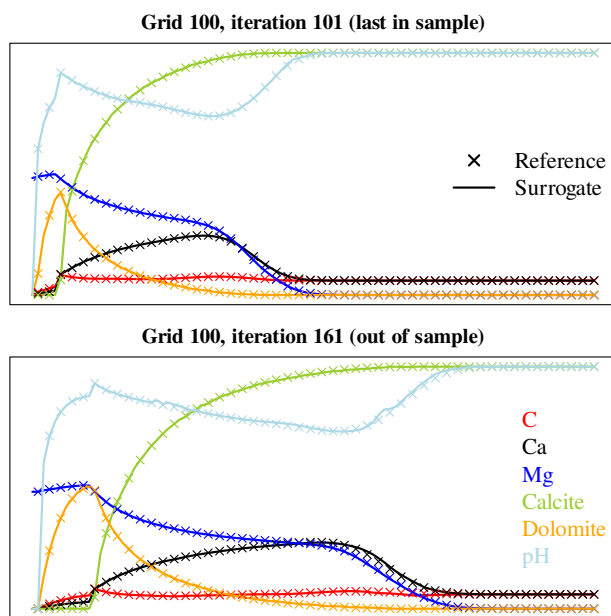


Figure 2. Profiles of total concentrations, pH and minerals for reference and hierarchical coupling 1D simulations with tolerance on mass balance error set to 10^{-6} , for grid 100. The top row displays the last *in sample* time step, the bottom the last simulated time step, after 60 iterations for which the surrogate was *out of sample*.

“out of sample”, in the sense of seen vs. unseen geochemical data, since the training data only comprises “physical” input combinations but, due to the introduced inaccuracies, we are asking the surrogate more and more predictions based on slightly “unphysical” input combinations. Having highly accurate surrogates, hence, would be beneficial also in this regard.

It is difficult to discriminate “a priori” between acceptable and unacceptable simulation results based on a threshold of an error measure such as that of eq. 3, which can be roughly interpreted as “mean percentage error”. This is also a point where in our opinion further research is needed. Relying on the visualisation of the surrogate simulation results and reference, we can summarize that the tolerance on mass balance of 10^{-6} (solid lines in Figure 3) produces accurate coupled simulations, excellent accuracy within the time steps of the training data and good accuracy after the 60 out of sample iterations. The tolerance of 10^{-5} as well as the simulations based solely on surrogates produce acceptable accuracy in sample but unusable and rapidly diverging results out of sample. For the given chemical problem, the 10^{-6} tolerance on mass balance could be relaxed, whereas the 10^{-5} is too optimistic. The optimal value, at least for the considered time steps, lies between these two values.

The overall speedup - in terms of total wall clock time of the coupled simulations, thus including also CPU time used for advection and all the overheads, although both much less computationally intensive than chemistry, and therefore termed pseudo speedup - with respect to the reference simulations is summarized in Figure 4. Here the whole 161 iterations, also all the out of sample ones are considered. Pseudo speedup increases with grid size as expected. The accurate 10^{-6} simulations are not accelerated on grid 50 (pseudo speedup of 0.86), but they reach 1.33 on the 200 grid. The surrogate-only speedup

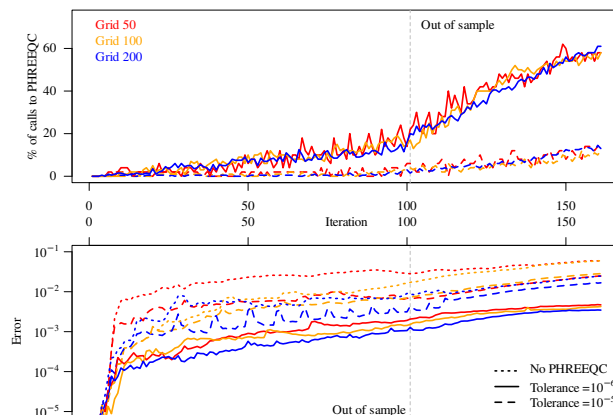


Figure 3. Purely data-driven approach: evaluation of calls to full physics simulator for the runs with hierarchical coupling for the three discretisations à 50, 100 and 200 elements, and of overall discrepancy between surrogate simulations and reference. When the surrogate enters the region of “unseen data”, its accuracy degrades significantly, which causes loss of efficiency rather than accuracy.

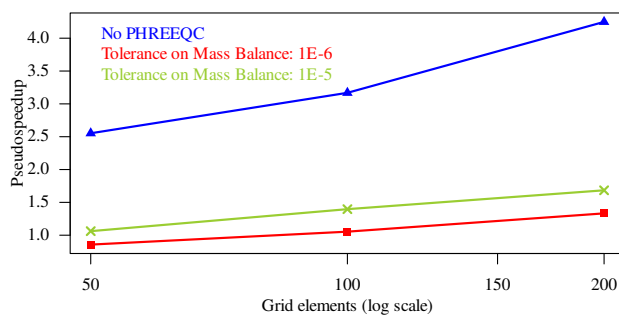


Figure 4. Overall pseudo speedup (total wall clock time) after 161 iterations for coupled simulations with hierarchical coupling and only relying on the surrogate.

starts at around 2.6 for the 50 grid and reaches 4.2 for the 200 grid, and that can be taken as a measure of achievable speedup, in projection, by using extremely accurate surrogates. Considering only the first 101 iterations, the 10^{-6} simulations would achieve speedup slightly larger than one already on the 50 elements grid, and be well over 2 on the 200 grid. Extrapolating to grids with 10^5 or 10^6 elements, speedups in the order of 25-50 are achievable for this chemical problem. The speedup will arguably increase even further in presence of more complex chemistry.



4 Surrogates based on geochemical knowledge

The above presented fully data-driven approach disregards any domain knowledge or known physical relationships between variables besides those which are picked up automatically by the multivariate algorithms operating on the input/outputs in the training data.

We start this approach by considering the actual “true” degrees of freedom for the geochemical problem, which is fully described by seven inputs and four outputs: Δ calcite, Δ dolomite, Cl and pH. This means that we will have to calculate back the changes in concentrations for C, Ca and Mg, risking a quicker propagation of errors if the reaction rates of the minerals are incorrectly predicted.

The reference simulations for this part are run with $\nu = 1$ and thus without numerical dispersion on four different grids: 50, 100, 200 and 500 elements respectively. This implies that the simulation on grid 500 has ten times more coupling iterations than the 50 grid, or in other terms, that the allowable time step in grid 500 is a tenth of that for grid 50.

A common way to facilitate the task of the regressors by “injecting” physical knowledge into the learning task of the algorithms is to perform *feature engineering*: this simply means computing new variables defined by non-linear functions of the original ones, which may give further insights regarding multivariate dependencies, hidden conditions or relevant subsets of the original data.

For any geochemical problem involving dissolution or precipitation of minerals, each mineral’s saturation ratio (SR) or its logarithm SI (Saturation Index) discriminates the direction of the reaction. If $SR > 1$ (and thus $SI > 0$) the mineral is oversaturated and precipitates; it is undersaturated and dissolves if $SR < 1$ ($SI < 0$); $SR = 1$ ($SI = 0$) implies local thermodynamic equilibrium. Writing the reaction of calcite dissolution:



The Law of Mass Action (LMA) relates, at equilibrium, the activities of the species present in the equation. We conventionally indicate activity with square brackets. For eq. 5, the LMA reads:

$$\begin{aligned} K_{\text{Cc}}^{\text{eq}} &= \frac{[\text{Ca}^{+2}]_{\text{eq}} \cdot [\text{HCO}_3^-]_{\text{eq}}}{[\text{H}^+]_{\text{eq}}} \\ &= \frac{\text{Ca}_{\text{eq}}^{+2} \cdot \text{HCO}_3^-_{\text{eq}}}{[\text{H}^+]_{\text{eq}}} \cdot \gamma_{\text{Ca}^{+2}} \gamma_{\text{HCO}_3^-} \end{aligned} \quad (6)$$

where γ stands for the activity coefficient of subscripted aqueous species. The solubility product $K_{\text{Cc}}^{\text{eq}}$ at equilibrium, tabulated in thermodynamic databases, is a function of temperature and pressure and defines the saturation ratio:

$$\text{SR}_{\text{Cc}} = \frac{1}{K_{\text{Cc}}^{\text{eq}}} \frac{\text{Ca}^{+2} \cdot \text{HCO}_3^-}{[\text{H}^+]} \cdot \gamma_{\text{Ca}^{+2}} \gamma_{\text{HCO}_3^-} \quad (7)$$

The estimation of the saturation ratio of equation 7 using the elemental concentrations available in our training data is the first, natural feature engineering tentative. Hereby, a few assumptions must be made.

Using total elemental concentrations as proxy for species activities implies neglecting the actual speciation to estimate the ion activity products, but also the difference between concentration and activity - “true” activity $[\text{H}^+]$ is known from the pH -



but also not relying on. For the chemical problem at hand, as it will be shown, it is a viable approximation, but it will not be in presence of strong gradients of ionic strength or in general for more complex or concentrated systems. An exception to this simplification is required for dissolved carbon due to the well known buffer. In this case, given that the whole model is at pH between 7 to 10, we may assume that two single species dominate the dissolved carbon speciation: CO_3^{-2} and HCO_3^- . The relationship between the activities of those two species is always kept at equilibrium in the PHREEQC models and thus, up to the “perturbation” due to transport, also in our dataset. This relationship is expressed by the reaction and the corresponding law of mass action written in eq. 8:



The closure equation, expressing the approximation of total carbon concentration as the sum of two species, gives us the second equation for the two unknowns:

$$C = \text{HCO}_3^- + \text{CO}_3^{-2} \quad (9)$$

Combining eq. 8 and 9 we get the estimation of dissolved bicarbonate (the wide tilde indicates that it is an estimation) from the variables total carbon and pH comprised in our dataset and an externally calculated thermodynamic constant:

$$\widetilde{\text{HCO}_3^-} := \frac{C \cdot [\text{H}^+]}{K_{\text{carb}}^{\text{eq}} + [\text{H}^+]} \quad (10)$$

Now we can approximate the theoretical calcite saturation ratio $\widetilde{\text{SR}}_{\text{Cc}}^{\text{theor}}$ with the formula:

$$\widetilde{\text{SR}}_{\text{Cc}}^{\text{theor}} := \frac{\text{Ca} \cdot \widetilde{\text{HCO}_3^-}}{[\text{H}^+] \cdot K_{\text{Cc}}^{\text{eq}}} = \frac{\text{Ca} \cdot C}{K_{\text{Cc}}^{\text{eq}} (K_{\text{carb}}^{\text{eq}} + [\text{H}^+])} \quad (11)$$

The two thermodynamic quantities (at 25 °C and atmospheric pressure) $K_{\text{carb}}^{\text{eq}} = 10^{-10.3288}$ and $K_{\text{Cc}}^{\text{eq}} = 10^{-2.00135}$ were computed with the CHNOSZ package for the R environment (Dick, 2019), but may also be derived with simple algebraic calculations from, i.e., the same PHREEQC database employed in the reactive transport simulations.

Do these two newly defined variables, or “engineered features”, bicarbonate and calcite saturation ratio, actually help to better understand and characterize our dataset? This can be simply assessed by plotting the Δ_{calcite} against the logarithm of $\widetilde{\text{SR}}_{\text{Cc}}^{\text{theor}}$, which is the $\widetilde{\text{SI}}_{\text{Cc}}^{\text{theor}}$ (Figure 5a, leftmost panel, dataset from the reference simulations on grid 200, which will be used from now on to illustrate the analysis since it contains enough data points) in the data. While many points remarkably lie on a smooth curve (colored in black), many others are scattered throughout the graph (in red). It is easy to observe that those red points are either on the trivial $\Delta_{\text{calcite}}=0$ line, implying that calcite is undersaturated but not present in the system so nothing happens, or else the reaction did not reach the amount which could have been expected based on its initial undersaturation simply because calcite has been completely depleted during the timestep. All the red points correspond in facts to simulations with calcite=0 in the labels (results) dataset. The retained black points, however, belong to timesteps where the dissolution of calcite is limited by kinetics and not by its initial amount, and can be thus used to estimate the reaction rate.



Figure 5a also displays a problem with the defined $\widetilde{SR}_{Cc}^{theor}$: its relationship is not bijective with the $\Delta_{calcite}$. This means that we should proceed now to split the data in two different regions above and under the cusp (signalled by the blue horizontal line). However, just simply dropping the denominator of equation 11 solves this problem to a large extent:

$$345 \quad \widetilde{SR}_{Cc} := Ca \cdot C \quad (12)$$

The center panel of Figure 5b shows the scatter plot of $\Delta_{calcite}$ versus the simplified \widetilde{SI}_{Cc} . All points lie now on a smooth curve, and the relation between the two variables is indeed quite perfectly bijective, with the exception of points very close to the $\Delta_{calcite}=0$ line, where they are more scattered; but since those points also correspond to the smallest amounts of reactions, we can deem this as successful approximation. Note that dropping the denominator in the definition of \widetilde{SR}_{Cc} also means that this feature does not reach one at equilibrium (and \widetilde{SI}_{Cc} zero), which is clear observing the range of the x-axis in panels a and b of Figure 5. This has however no practical consequence for this problem: calcite is always undersaturated or at equilibrium in the benchmark, and we just defined a simple feature which is in bijective relationship with the amount of “true” dissolution in the data. While it could be possible to derive an analytical functional dependency between the observed amount of dissolved calcite and the estimated \widetilde{SI}_{Cc} , for example manipulating the kinetic law, we opted to use a regressor instead. The good bijectivity between the two variables means that we should be able to regress the first using only the second. In the rightmost panel of Figure 5c are plotted in blue the *in sample* predictions of a Multivariate Adaptive Regression Spline model (MARS) (Friedman, 1991, 1993), computed through the `earth` R package (Milborrow, 2018), based only on \widetilde{SI}_{Cc} . The accuracy is already acceptable indeed; however including further predictors from the already available features, in this case pH, Ca and Mg, a better regression (in red) is achieved, improving the RMSE of more than factor two.

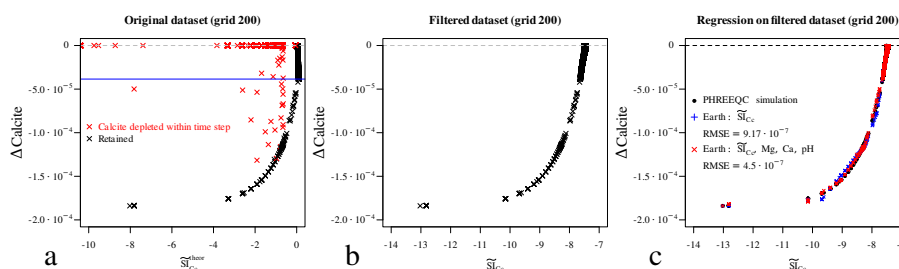


Figure 5. (a) Scatter plot of $\Delta_{calcite}$ vs estimated $\widetilde{SI}_{Cc}^{theor}$. The datapoints in red cannot be used to estimate the reaction rate from the dataset since calcite is depleted within the simulation time step. Furthermore, the retained points are not in bijective relationship with the $\Delta_{calcite}$, with the blue horizontal line separating two regions where bijectivity is given. (b) $\Delta_{calcite}$ versus the simplified \widetilde{SI}_{Cc} : bijectivity is achieved. (c) A MARS regressor is computed for the retained data points, based solely on the estimated \widetilde{SI}_{Cc} (in blue) and using also other predictors to ameliorate the multivariate regression.

360 Before moving forward, two considerations are important. First, the red points of Figure 5a should not be used when trying to estimate the rate of calcite dissolution, since they result from a steep and “hidden” non-linearity or discontinuity in the



underlying model. This is a typical example of data potentially leading to overfitting in a machine-learning sense. Secondly, this “filter” does not need to be applied at runtime during coupled reactive transport simulations: it suffices to estimate correctly the reaction rate given the initial state and then ensure that calcite does not reach negative values.

365 More interesting and more demanding is the case of dolomite, which firstly precipitates and then re-dissolves in the benchmark simulations. In a completely analogous manner as above we define its saturation ratio \widetilde{SR}_{Dol} as:

$$\widetilde{SR}_{Dol} = \frac{Mg \cdot Ca \cdot C^2}{[H^+]^2 \cdot K_{Dol}^{eq}} \quad (13)$$

thus using the total elemental dissolved concentration of C and with $K_{Dol}^{eq} = 10^{3.647}$ resulting from the reaction:



370 The theoretical value of $K_{Dol}^{eq} = 10^{3.647}$ used for calculation of \widetilde{SI}_{Dol} does not discriminate the initially undersaturated from the oversaturated samples (dashed vertical black line in Figure 6). The “offset” which would serve us for a correct discrimination is nothing else than the maximum value of \widetilde{SR}_{Dol} restricted to the region where $\Delta_{dolomite} \leq 0$. We correspondingly update the definition of \widetilde{SR}_{Dol} :

$$\widetilde{SR}_{Dol} = \frac{Mg \cdot Ca \cdot C^2}{[H^+]^2 \cdot K_{Dol}^{eq}} - \max(\widetilde{SR}_{Dol} |_{\Delta_{dolomite} \leq 0}) \quad (15)$$

Now we are guaranteed that the vertical line $\widetilde{SR}_{Dol}=1$ (or equivalently, $\widetilde{SI}_{Dol}=0$, plotted with a dashed blue line in Figure 6)

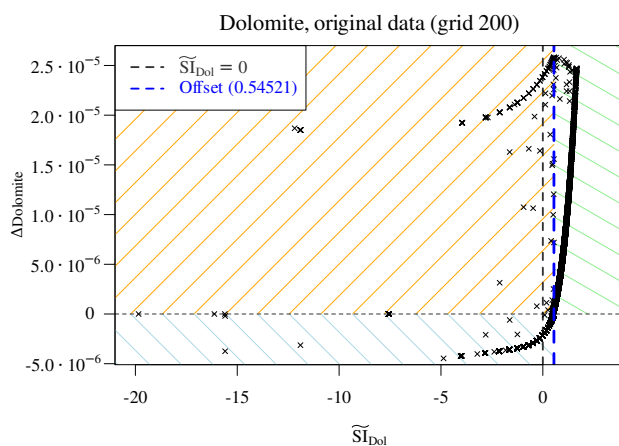


Figure 6. Scatter plot of $\Delta_{dolomite}$ vs estimated \widetilde{SI}_{Dol} . The theoretical $\widetilde{SI}_{Dol}=0$ does not discriminate the initially undersaturated from the oversaturated samples (dashed vertical black line), and must be corrected with an apparent offset (blue dashed line). The plot identifies three distinct regions in parameter space: initially supersaturated and precipitating dolomite (top right, green shading); initially undersaturated and dissolving (bottom left, blue shading); and points where dolomite is initially undersaturated but ends up precipitating (top left, orange shading).

375

divides correctly the parameter space in four distinct quadrants. Note that this offset emerges from the actual considered data,



and depends on the perturbation of the concentrations due to transport and thus, in our simple advective scheme, on the grid resolution through the time step. It follows that a different offset is expected for the other grids, and a different learning for each grid is necessary.

380 The green shaded, top right quadrant points to dolomite precipitation in initially supersaturated samples; the bottom left, blue shaded contains solutions initially undersaturated w.r.t. dolomite and, if present, dissolving; the top left, orange shaded quadrant is the most problematic: dolomite is initially undersaturated but, presumably due to the concurring dissolution of calcite, it becomes supersaturated during the time step and hence precipitates.

First of all, we note that the initial presence of calcite is a perfect proxy for $\tilde{S}I_{Dol}$. If calcite is initially present in points
 385 reached by the reactive magnesium chloride solution, then dolomite precipitates. When calcite is completely depleted, then dolomite starts dissolving again. The dissolution of dolomite in absence of calcite follows the same logic as the dissolution of calcite above: a few points are scattered inbetween the line $\Delta dolomite=0$ and the envelope of points lying on a well defined curve. These scattered points are again those where dolomite is depleted within the time step, so they are excluded. For the remaining points, an xgboost regressor based on the predictors $\tilde{S}I_{Dol}$, pH, C, Cl, Mg and dolomite achieves an excellent accuracy (Figure 7) in reproducing the observed $\Delta dolomite$. The top right quadrant of Figure 6, corresponding to the case of

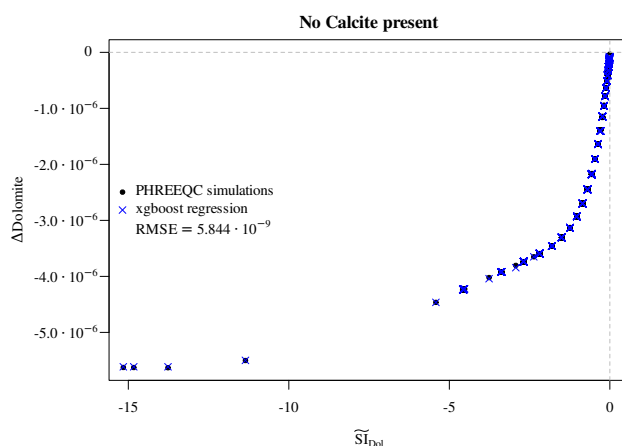


Figure 7. Regression of $\Delta dolomite$ vs estimated $\tilde{S}I_{Dol}$ for the cases where no calcite is initially present. The multivariate regressor makes use of the predictors $\tilde{S}R_{Dol}$, pH, C, Cl, Mg and dolomite.

390 dolomite precipitating while calcite is dissolving, cannot be explained based only on the estimated $\tilde{S}I_{Dol}$ since their relationship is not surjective (Figure 8a). Here again we can use a piece of domain knowledge to engineer a new feature to move forward. The Mg/Ca ratio is often used to study the thermodynamics of dissolution of calcite and precipitation of dolomite (Möller and De Lucia, 2020). Effectively, the occurring overall reaction which transforms calcite into dolomite reads:



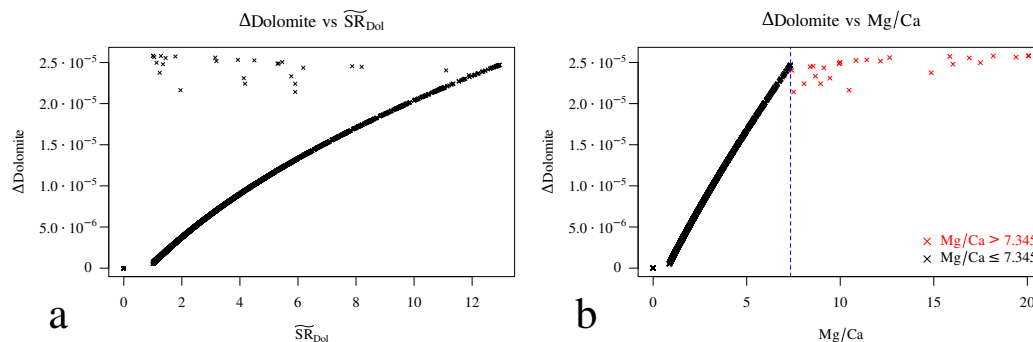


Figure 8. Precipitation of dolomite in presence of calcite. (a) the relationship between Δ dolomite and its saturation ratio is not surjective. (b) The Mg/Ca ratio perfectly discriminates two distinct regions in parameter space.

By applying the law of mass action to reaction 16, it is apparent that its equilibrium constant is a function of the Mg/Ca ratio (of its inverse in the form of equation 16). Plotting the Δ dolomite versus the initial Mg/Ca ratio, a particular ratio of 7.345 discriminates between two distinct regions for this reaction. Incidentally, this splitting value corresponds to the highest observed \widetilde{SR}_{Dol} in the training data; again, as previously noted for the offset on the estimated saturation index, this numerical value depends on the considered grid and time step. On the left hand region we observe a smooth, quasi-linear dependency of the amount of precipitated dolomite on initial Mg/Ca. This is a simple bijective relationship where we can apply a simple monivariate regression. The amount of precipitated dolomite is accurately predicted by a MARS regressor using the sole Mg/Ca as predictor.

The region on the right of the splitting ratio can be best understood considering the fact that the precipitation of dolomite is limited, in this region, by a concurrent amount of calcite dissolution. The full physics chemical solver finds iteratively the correct amounts of calcite dissolution and dolomite precipitation while honoring both the kinetic laws and all the other conditions for a geochemical DAE system (mass action equations, electroneutrality, positive concentrations, activity coefficients, ...). We cannot reproduce such articulate and “interdependent” behavior without knowing the actual amount of dissolved calcite: we are forced here to employ the previously estimated Δ calcite as a “new feature” to estimate of the amount of dolomite precipitation, albeit limited to this particular region of the parameter space. A surprisingly simple expression, fortunately, captures this relationship quite accurately (Figure 9). This implies of course that during coupled simulations first the Δ calcite must be computed, and relying on this value, the Δ dolomite can be further estimated.

The last parameter space region which is left to consider is the orange-shaded, topleft quadrant of Figure 6. Here, although dolomite is undersaturated at the beginning of the time step, it still precipitates in the end, following the concurrent dissolution of calcite which changes its saturation state. Since however we already calculated the Δ calcite, we can update the concentrations of dissolved Ca and C of corresponding amounts. One of these two concentrations, together with that of Mg, will constitute a limiting factor for the precipitation of dolomite. Hence, plotting the Δ dolomite against the minimum value of these three concentrations at each point (C must be divided by two for the stoichiometry of dolomite), we obtain a piecewise-

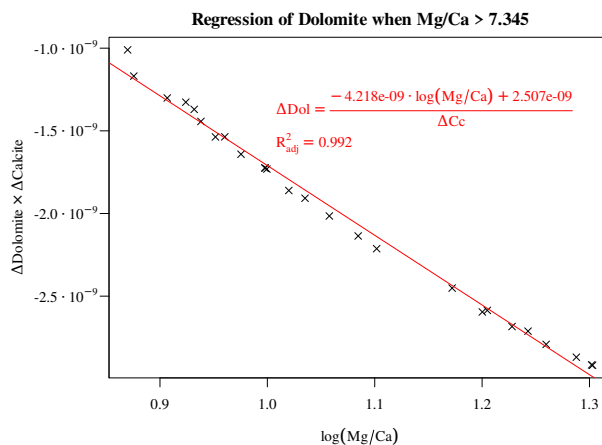


Figure 9. Regression of Δ dolomite in the right hand region of Figure 8b.

linear relationship with limited non-linear effects. A very simple regression is hence sufficient to capture the bulk of the “true model behavior” for all these data points (Figure 10). Now the behavior of calcite and dolomite is fully understood and we dis-

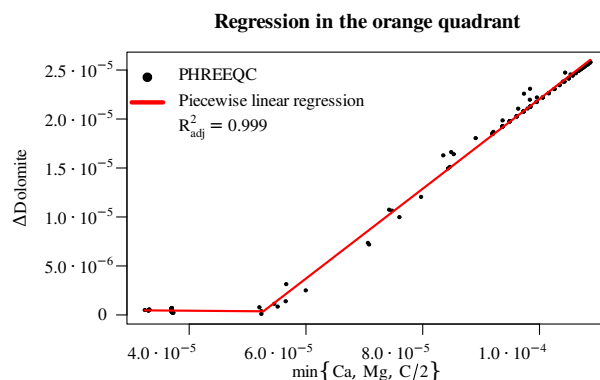


Figure 10. Piecewise-linear regression for the orange-shaded, top left quadrant of figure 6 based on the limiting elemental concentration after having considered calcite dissolution.

420

pose of a surrogate for both of them. Among the remaining output variables, only pH needs to be regressed: Cl is non-reactive, meaning that the surrogate is the identity function. For pH, while it could be possible to derive a simplified regression informed with geochemical knowledge, we chose for simplicity to use the xgboost regressor.

Summarizing, we effectively designed a *decision tree*, based on domain knowledge, which enabled us to make sense of the “true” data, to perform physically meaningful feature engineering and ultimately to define a surrogate model “translated” to the data domain (Figure 11). The training of this decision tree surrogate consists merely in computing the engineered features,

425

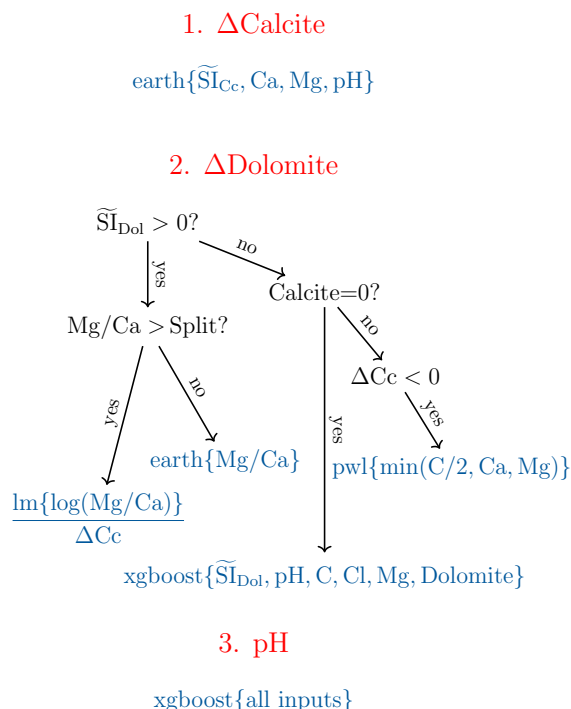


Figure 11. Decision tree for the surrogate based on physical interpretation of the training dataset. The engineered features are used as splits and as predictors for different regressions depending on the region of parameter space. The abbreviations “lm” and “pwl” stand respectively for “linear model” and for piecewise-linear regression.

finding the apparent offset for the $\widetilde{\text{SI}}_{\text{Dol}}$, the split value for the Mg/Ca ratio, and performing six distinct regressions on data subsets, of which three are monivariate and two use less predictors than the corresponding completely data-driven counterpart. All of them, excluding pH, only use a subset of the original training dataset. On our workstation, this operation takes few
 430 seconds. The resulting surrogate is valid for the Δt of the corresponding training data.

To evaluate the performance of this surrogate approach, a decision tree is trained separately for each grid (and hence Δt) using the reference timesteps until 42000 seconds, whereas the coupled simulations are prolonged to 60000 s, so that at least 30 % of the simulation time is computed on unseen geochemical data.

The top panel of Figure 12 shows the results of the coupled simulation for grid 50 using the surrogate trained on the same
 435 data, at the end of the iterations used for training. Discrepancies with respect to the reference full physics simulation are already evident. The problem here is that the training dataset is too small and the time step too large for the decision tree surrogate to be accurate. However, nothing forbids to perform “inner iterations” for the chemistry using a surrogate trained on a finer grid, which directly corresponds to smaller Δt . For grid 50 ($\Delta t=1066$ s) we can hence use the surrogate trained on grid 500 ($\Delta t=106.6$ s) just calling it 10 times within each coupling iterations. The bottom panel of Figure 12 displays the corresponding



results. The same problem affects the grid 100, which also requires the surrogate trained on grid 500, reiterated 5 times in this

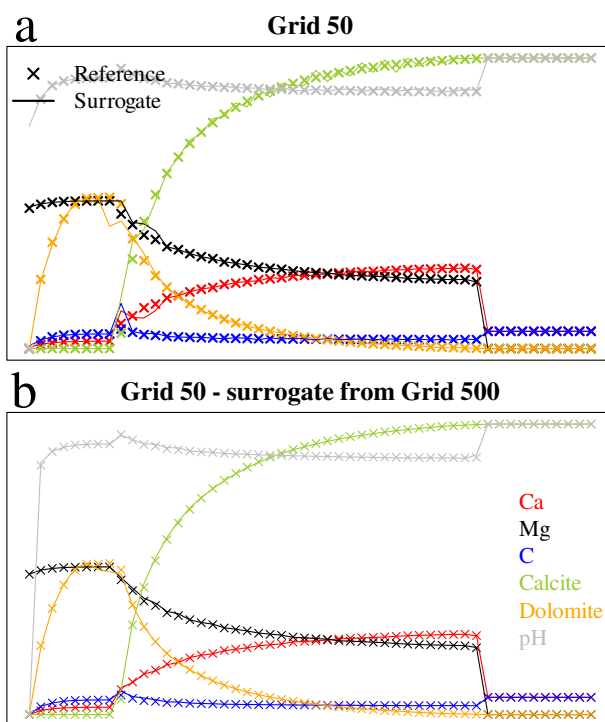


Figure 12. Comparison of variables profiles for coupled simulations using the decision tree approach versus the references, at the end of the timesteps used for training for grid 50 (41 coupled iterations). (a) decision tree trained on the data from reference grid 50 ($\Delta t=1066$ s). (b) surrogate simulations using decision tree trained on grid 500 ($\Delta t=106.6$ s), repeated 10 times for each coupling time step.

440

case. The grids 200 and 500 are fine with the own reference data, as can be seen in Figure 13, this time displaying the end of simulation time at 60000 s.

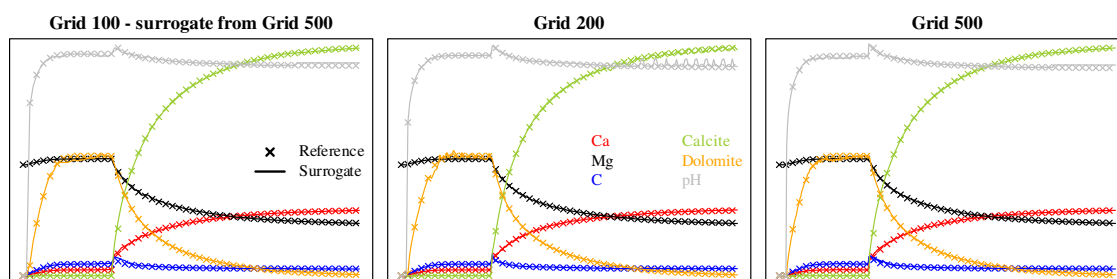


Figure 13. Variable profiles after 60000 s (simulation time) for grids 100, 200 and 500.



In Figure 14 are summarised the errors of the surrogates simulations (top panel) and the overall pseudo speedup after 60000 s (bottom panel). While inaccuracies are indeed introduced in the coupled simulations by the decision tree surrogate, crossing the “out of sample” boundary does not provoke a steep increase in error. Even if the overall error is slightly larger than the corresponding purely data-driven simulations with 10^{-6} tolerance, the physics-based approach has the major advantage of being much more robust when encountering unseen data. Moreover, since no calls to PHREEQC are issued at all during these simulations, the performance of the coupled simulations is not going to degrade during the simulation time. The physics-based surrogates achieve large pseudo speedups, starting with 2.7 for the grid 50 and reaching 6.8 for the 500 grid (Figure 14, bottom panel).

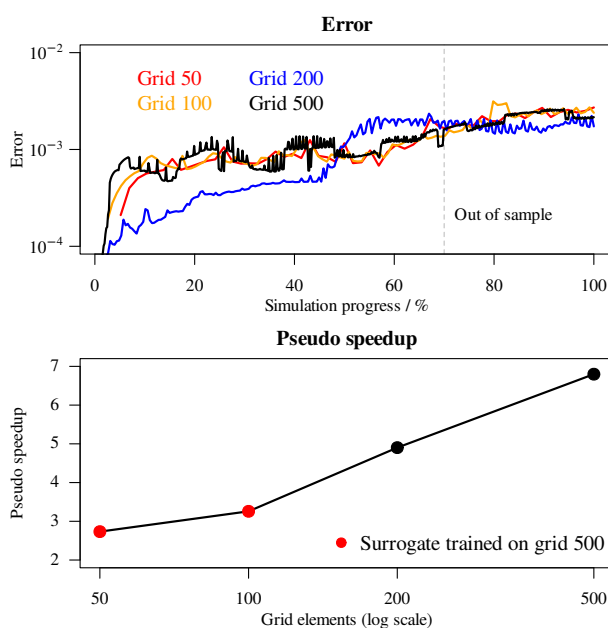


Figure 14. Top: errors of surrogate simulations w.r.t. references. Bottom: overall pseudo speedup after 60000 s.

Note that the decision tree approach has been implemented in pure high-level R language (up to the calls to the regressors `xgboost` and `earth`, which are implemented in low-level languages such as C/C++) and is not optimized. A better implementation would further improve its performance, especially in the case where repeated calls to the surrogate are performed at each coupled iteration.

455 5 Discussion and future work

The results presented in this work devise some strategies which can be successfully exploited to speedup reactive transport simulations. The simplifications concerning the transport and the coupling itself (stationary flow; pure advection with dispersive full explicit forward Euler scheme; no feedback of chemistry on porosity and permeability; initially homogeneous medium;



kinetic rate not depending on reactive surfaces) are obviously severe, but most of them should only marginally affect the
460 validity of the benchmarks concerning the achievable speedup of geochemistry in a broad class of problems.

A fully data-driven approach, combined with a hierarchical coupling in which full physics simulations are performed only if
surrogate predictions are found implausible, is feasible and promises significant speedups for large scale problems. The main
advantage of this approach is that the same “code infrastructure” can be used to replace any physical process, not limited
to geochemistry: it is completely general, and it could be implemented in any multiphysics toolbox to be used for any co-
465 simulated process. The hierarchy of models for process co-simulation is a vast research field on itself. This idea has to our
knowledge never been implemented specifically for reactive transport, but has been proposed, e.g., for particular problem
settings in fluid dynamics and elastomechanics (Altmann, 2013; Altmann and Heiland, 2015) and in the broader context of
theoretical model reduction and error control (Domschke et al., 2011). This is however a fertile interdisciplinary research task
and it is not difficult to foresee that significant progress in this area will soon be required to facilitate and fully leverage the
470 powerful machine learning algorithms already available, in order to speedup any complex, multiscale numerical simulations.

The coupling hierarchy implemented in this work is obviously extremely simple and cannot be directly compared with the
above cited works, since it is merely based on a *a posteriori* evaluation of plausibility of geochemical simulations. Furthermore, it
exploits redundant regressions, which is suboptimal, albeit practical: in effects, regressing more variables than strictly necessary
is not much different than regressing the true independent variables and their error models. Since the surrogate predictions are
475 so cheap compared to the full physics, it would be only slightly beneficial to first interrogate the error model and then go
directly to the full physics instead of computing at once the whole surrogate predictions and check it afterwards. Nevertheless,
several improvements can be implemented with respect to the hierarchy presented in this work. The first would be to add charge
balance to the error check at runtime. For different classes of chemical processes, other criteria may be required. For example
check on mass action laws can be implemented for models requiring explicit speciation, like in the simulations of radionuclide
480 diffusion and sorption in storage formations. Another one would be to actually eliminate one or more redundant regression
and base the error check on the accordance between the overlapping one. As an example, one could regress the Δ dolomite,
 Δ calcite and Δ Ca, limiting in practice the mass balance check to one element.

In our opinion there is no point in discussing if there is one most suitable or most efficient regression algorithm. This largely
depends on the problem at hand and on the skills of the modeller. While we rather focused on gradient boosting decision-tree
485 regressors for the reasons briefly discussed in section 3, a consistent number of authors successfully applied artificial neural
networks to a variety of geochemical problems and coupled simulations (Laloy and Jacques, 2019; Guérillot and Bruyelle,
2020; Prasianakis et al., 2020). We would like to point out that transforming geochemistry - as any other process - in a
pure machine learning problem requires on one hand skills that are usually difficult for the geoscientists to acquire, and on
the other it fatally overlooks domain knowledge that can be used to improve at least the learning task, which will directly
490 result in accurate and robust predictions, as we demonstrated in section 4. Feature engineering based on known physical
relationships and equations should be part of any machine learning workflow anyway; building experience in this matter,
devising suitable strategies for a broad class of geochemical problems is in our opinion much more profitable than trying
to tune overly complex “black box” models of general applicability. The purely data-driven approach has its own rights and



495 applications. As already noted, it is a completely process-agnostic approach which can be implemented in any simulator for
any process. However, in absence of physical knowledge *within* the surrogate, the training data must cover beforehand all the
processes and the scenarios happening in the coupled simulations. On-demand training and successive incremental update of
the surrogates at runtime during the coupled simulations would mitigate this issue. This would require a careful choice of
the regressors, since not all of them have this capability, and possibly a sophisticated load balance distribution, likely viable
only in the context of parallel computing. In perspective, however, this is a feature that in our opinion should be implemented
500 in the numerical simulators. A second issue, related to the first, is that a data-driven surrogate trained on a specific chemical
problem (intending here initial conditions, concentration of the injected solutions, mineral abundances, time steps...), is not
automatically transferable to different problem settings, even when for example only a single kinetic constant is varied. Again,
shaping the surrogate following the physical process to be simulated seem here the most straightforward way to overcome this
issue, at least partially. One would dispose of partial regressions in specific parameter space regions which could be varied
505 following changes in underlying parameters.

It remains to be assessed if it is possible to generalize and automate the physics-based surrogate approach devised in sec-
tion 4 on geochemical problems of higher complexity, i.e., with many minerals reacting. No claim of optimality is made about
the actual choice of engineered features we made for this chemical benchmark: different features could possibly explain the
data even more simply, and thus the chemical process. The important part is the principle: identify relationships as bijective
510 as possible between input and output parameters, compartmentalized in separated regions of parameter space, using features
derived by the governing equations. An automation of feature engineering based on stoichiometry of the minerals is a straight-
forward extension, since it can be achieved by simply parsing the thermodynamic databases. An automatic application of the
approach starting with a large number of engineered features may originate forests of trees much like the well known random
forest or gradient boosting algorithms, but specialised in geochemical models: a true hybrid physics-AI model.

515 Also, the regressors which constitute the leaves of the decision tree of Figure 11 are completely arbitrary and were selected
based on our own experience. A more in-depth breakdown of the relationships between variables, for example analytical
expressions derived directly from the kinetic law, could reduce most or all regressions to simple statistical linear models, which
would even further increase the interpretability of the surrogate.

6 Conclusions

520 Employing surrogates to replace computationally intensive geochemical calculations is a viable strategy to speedup reactive
transport simulations. A hierarchical coupling of geochemical sub-processes, allowing to recur to “full physics” simulations
when surrogate predictions are not accurate enough is advantageous to mitigate the inevitable inaccuracies introduced by the
approximated surrogate solutions. In the case of purely data-driven surrogates, which are a completely general approach not
limited to geochemistry, regressors operate exclusively on input/output data oblivious to known relationships. Here, redundant
525 information content can be employed effectively to obtain cheap estimation of plausibility of surrogate predictions at runtime,
by checking the errors on mass balance. This estimation works well at least for the presented geochemical benchmark. Our



tests show consistent advantage of decision-tree based regression algorithms, especially belonging to the gradient boosting family.

530 Feature engineering based on domain knowledge, i.e., the actual governing equations for the chemical problem as solved by the full physics simulator, can be used to construct a surrogate approach in which the learning task is enormously reduced. The strategy consists in partitioning the parameter space based on the engineered features, and looking for bijective relationships within each region. This approach reduces both the required multivariate predictions and the dimension of training dataset upon which each regressor must operate. Algorithmically it can be represented by a decision tree, and proved both accurate and robust, being equipped to handle unseen data and less sensible to sparse training dataset, since it embeds and exploits
535 knowledge about the modelled process. Further research is required in order to generalize it and to automate it for more complex chemical problems, as well as for adapting it to specific needs such as sensitivity and uncertainty analysis.

Both approaches constitute non-mutually-exclusive valid strategies in the arsenal of modellers dealing with overwhelmingly CPU-expensive reactive transport simulations, required by present day challenges in subsurface utilisation. In particular, we are persuaded that hybrid AI-physics models will offer the decisive computational advantage needed to overcome current
540 limitations of classical equation-based numerical modelling.

Code availability. DecTree v1.0 is a model experiment setup and evaluated in the R environment. All the code used in the present work is available under LGPL v2.1 license at Zenodo with DOI 10.5281/zenodo.4569573 (<https://doi.org/10.5281/zenodo.4569573>)

DecTree depends on the RedModRphree package v0.0.4, equally available at Zenodo with DOI 10.5281/zenodo.4569516 (<https://doi.org/10.5281/zenodo.4569516>).

545 *Author contributions.* MDL shaped the research, performed analyses, programming and wrote the manuscript. MK helped providing funding, shaping the research, and revised the manuscript.

Competing interests. The authors have no conflict of interests.

Acknowledgements. The authors gratefully acknowledge the Helmholtz Association of German Research Centers - Initiative and Networking Fund for the funding in the framework of the project “*Reduced Complexity Models* – Explore advanced data science techniques to create
550 models of reduced complexity” - reference number ZT-I-0010.



References

- Altmann, R.: Index reduction for operator differential-algebraic equations in elastodynamics, *Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 93, 648–664, <https://doi.org/10.1002/zamm.201200125>, 2013.
- Altmann, R. and Heiland, J.: Finite element decomposition and minimal extension for flow equations, *ESAIM Math. Model. Numer. Anal.*, 49, 1489–1509, <https://doi.org/10.1051/m2an/2015029>, 2015.
- Appelo, C. A. J., Parkhurst, D. L., and Post, V. E. A.: Equations for calculating hydrogeochemical reactions of minerals and gases such as CO₂ at high pressures and temperatures, *Geochimica et Cosmochimica Acta*, 125, 49 – 67, <https://doi.org/10.1016/j.gca.2013.10.003>, 2013.
- Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, <https://doi.org/10.1145/2939672.2939785>, 2016.
- De Lucia, M. and Kühn, M.: Coupling R and PHREEQC: Efficient Programming of Geochemical Models, *Energy Procedia*, 40, 464 – 471, <https://doi.org/10.1016/j.egypro.2013.08.053>, 2013.
- De Lucia, M., Lagneau, V., Fouquet, C. d., and Bruno, R.: The influence of spatial variability on 2D reactive transport simulations, *Comptes Rendus Geoscience*, 343, 406 – 416, <https://doi.org/10.1016/j.crte.2011.04.003>, 2011.
- De Lucia, M., Kempka, T., and Kühn, M.: A coupling alternative to reactive transport simulations for long-term prediction of chemical reactions in heterogeneous CO₂ storage systems, *Geoscientific Model Development*, 8, 279–294, <https://doi.org/10.5194/gmd-8-279-2015>, 2015.
- De Lucia, M., Kempka, T., Jatnieks, J., and Kühn, M.: Integrating surrogate models into subsurface simulation framework allows computation of complex reactive transport scenarios, *Energy Procedia*, 125, 580–587, <https://doi.org/10.1016/j.egypro.2017.08.200>, 2017.
- Dick, J. M.: CHNOSZ: Thermodynamic Calculations and Diagrams for Geochemistry, *Frontiers in Earth Science*, 7, <https://doi.org/10.3389/feart.2019.00180>, 2019.
- Domschke, P., Kolb, O., and Lang, J.: Adjoint-Based Control of Model and Discretization Errors for Gas Flow in Networks, *International Journal of Mathematical Modelling and Numerical Optimisation*, 2, 175–193, <https://doi.org/10.1504/IJMMNO.2011.039427>, 2011.
- Engesgaard, P. and Kipp, K. L.: A geochemical transport model for redox-controlled movement of mineral fronts in groundwater flow systems: A case of nitrate removal by oxidation of pyrite, *Water Resources Research*, 28, 2829–2843, <https://doi.org/10.1029/92WR01264>, 1992.
- Friedman, J. H.: *Multivariate Adaptive Regression Splines (with discussion)*, *Annals of Statistics* 19/1, Stanford University, <https://statistics.stanford.edu/research/multivariate-adaptive-regression-splines>, 1991.
- Friedman, J. H.: *Multivariate Adaptive Regression Splines (with discussion)*, Technical Report 110, Stanford University, Department of Statistics, <https://statistics.stanford.edu/research/fast-mars>, 1993.
- Guérillot, D. and Bruyelle, J.: Geochemical equilibrium determination using an artificial neural network in compositional reservoir flow simulation, *Computational Geosciences*, 24, 697–707, <https://doi.org/10.1007/s10596-019-09861-4>, 2020.
- Hassine, A., Masmoudi, A., and Ghribi, A.: Tweedie regression model: a proposed statistical approach for modelling indoor signal path loss, *International Journal of Numerical Modelling: Electronic Networks, Devices and Fields*, 30, e2243, <https://doi.org/10.1002/jnm.2243>, 2017.
- Jatnieks, J., De Lucia, M., Dransch, D., and Sips, M.: Data-driven Surrogate Model Approach for Improving the Performance of Reactive Transport Simulations, *Energy Procedia*, 97, 447–453, <https://doi.org/10.1016/j.egypro.2016.10.047>, 2016.



- Jørgensen, B.: Exponential Dispersion Models, *Journal of the Royal Statistical Society. Series B (Methodological)*, 49, 127–162, <https://doi.org/10.2307/2345415>, <http://www.jstor.org/stable/2345415>, 1987.
- 590 Laloy, E. and Jacques, D.: Emulation of CPU-demanding reactive transport models: a comparison of Gaussian processes, polynomial chaos expansion, and deep neural networks, *Computational Geosciences*, 23, 1193–1215, <https://doi.org/10.1007/s10596-019-09875-y>, 2019.
- Leal, A. M. M., Kyas, S., Kulik, D. A., and Saar, M. O.: Accelerating Reactive Transport Modeling: On-Demand Machine Learning Algorithm for Chemical Equilibrium Calculations, *Transport in Porous Media*, 133, 161–204, <https://doi.org/10.1007/s11242-020-01412-1>, 2020.
- 595 Milborrow, S.: earth: Multivariate Adaptive Regression Splines derived from mda::mars by T. Hastie and R. Tibshirani, <https://CRAN.R-project.org/package=earth>, r package, 2018.
- Möller, P. and De Lucia, M.: The impact of Mg²⁺ ions on equilibration of Mg-Ca carbonates in groundwater and brines, *Geochemistry*, 80, 125 611, <https://doi.org/10.1016/j.chemer.2020.125611>, 2020.
- Palandri, J. L. and Kharaka, Y. K.: A compilation of rate parameters of water-mineral interaction kinetics for application to geochemical modeling, Tech. rep., USGS Menlo Park, California, USA, 2004.
- 600 Parkhurst, D. L. and Wissmeier, L.: PhreeqcRM: A reaction module for transport simulators based on the geochemical model PHREEQC, *Advances in Water Resources*, 83, 176–189, <https://doi.org/10.1016/j.advwatres.2015.06.001>, 2015.
- Prasianakis, N. I., Haller, R., Mahrous, M., Poonosamy, J., Pfingsten, W., and Churakov, S. V.: Neural network based process coupling and parameter upscaling in reactive transport simulations, *Geochimica et Cosmochimica Acta*, <https://doi.org/10.1016/j.gca.2020.07.019>, 2020.
- 605 Prommer, H., Sun, J., and Kocar, B. D.: Using Reactive Transport Models to Quantify and Predict Groundwater Quality, *Elements*, 15, 87–92, <https://doi.org/10.2138/gselements.15.2.87>, 2019.
- R Core Team: R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>, 2020.
- 610 Shao, H., Dmytrieva, S. V., Kolditz, O., Kulik, D. A., Pfingsten, W., and Kosakowski, G.: Modeling reactive transport in non-ideal aqueous–solid solution system, *Applied Geochemistry*, 24, 1287–1300, <https://doi.org/10.1016/j.apgeochem.2009.04.001>, 2009.
- Steeffel, C. I., DePaolo, D. J., and Lichtner, P. C.: Reactive transport modeling: An essential tool and a new research approach for the Earth sciences, *Earth and Planetary Science Letters*, 240, 539–558, <https://doi.org/10.1016/j.epsl.2005.09.017>, 2005.
- Steeffel, C. I., Appelo, C. A. J., Arora, B., Jacques, D., Kalbacher, T., Kolditz, O., Lagneau, V., Lichtner, P. C., Mayer, K. U., Meeussen, J. C. L., Molins, S., Moulton, D., Shao, H., Šimůnek, J., Spycher, N., Yabusaki, S. B., and Yeh, G. T.: Reactive transport codes for subsurface environmental simulation, *Computational Geosciences*, 19, 445–478, <https://doi.org/10.1007/s10596-014-9443-x>, 2015.
- 615 Tweedie, M. C. K.: An index which distinguishes between some important exponential families. *Statistics: Applications and New Directions*, Proceedings of the Indian Statistical Institute, Golden Jubilee International Conference. Golden Jubilee International Conference (Eds. J. K. Ghosh and J. Roy). Calcutta: Indian Statistical Institute., *Statistics: Applications and New Directions.*, 579–604, 1984.