

GEMSFITS: Code package for optimization of geochemical model parameters and inverse modeling

Journal Article**Author(s):**

Miron, George D.; Kulik, Dmitrii A.; Dmytrieva, Svitlana V.; Wagner, Thomas

Publication date:

2015-04

Permanent link:

<https://doi.org/10.3929/ethz-b-000099508>

Rights / license:

[Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International](#)

Originally published in:

Applied Geochemistry 55, <https://doi.org/10.1016/j.apgeochem.2014.10.013>

This is the Green Open Access version of: Miron, G. D., Kulik, D. A., Dmytrieva, S. V., Wagner, T., 2013. GEMSFITS: Code package for optimization of geochemical model parameters and inverse modeling. *Applied Geochemistry*, v. 55, p. 28-45.
Original publication see: <https://doi.org/10.1016/j.apgeochem.2014.10.013>

GEMSFITS: Code package for optimization of geochemical model parameters and inverse modeling

George D. Miron^a, Dmitrii A. Kulik^b, Svitlana V. Dmytrieva^c, Thomas Wagner^d

^aInstitute of Geochemistry and Petrology, ETH Zurich, Switzerland, dan.miron@erdw.ethz.ch

^bLaboratory for Waste Management, Paul Scherrer Institut, 5232 Villigen PSI, Switzerland

^cInstitute of Environmental Geochemistry, Kyiv, Ukraine

^dDepartment of Geosciences and Geography, University of Helsinki, Finland

Revised version

Submitted to *Applied Geochemistry*

Date: 13 October 2014

Abstract

GEMSFITS is a new code package for fitting internally consistent input parameters of GEM (Gibbs energy minimization) geochemical-thermodynamic models against various types of experimental or geochemical data, and for performing inverse modeling tasks. It consists of the `gemsfit2` (parameter optimizer) and `gfshell2` (graphical user interface) programs both accessing a NoSQL database, all developed with flexibility, generality, efficiency, and user friendliness in mind. The parameter optimizer `gemsfit2` includes the GEMS3K chemical speciation solver (<http://gems.web.psi.ch/GEMS3K>), which features a comprehensive suite of non-ideal activity- and equation-of-state models of solution phases (aqueous electrolyte, gas and fluid mixtures, solid-solutions, (ad)sorption). The `gemsfit2` code uses the robust open-source NLOpt library for parameter fitting, which provides a selection between several nonlinear optimization algorithms (global, local, gradient-based), and supports a large-scale parallelization. The `gemsfit2` code can also perform comprehensive statistical analysis of the fitted parameters (basic statistics, sensitivity, Monte Carlo confidence intervals), thus supporting the user with powerful tools for evaluating the quality of the fits and the physical significance of the model parameters. The `gfshell2` code provides menu-driven setup of optimization options (data selection, properties to fit and their constraints, measured properties to compare with computed counterparts, and statistics). The practical utility, efficiency, and geochemical relevance of GEMSFITS is demonstrated by examples of typical classes of problems that include fitting of parameters of thermodynamic mixing models, optimization of standard state Gibbs energies of aqueous species and solid-solution end-members, thermobarometry, inverse titrations, and optimization problems that combine several parameter- and property types.

Keywords: parameter optimization, regression tool, thermodynamic modeling, Gibbs energy minimization, experimental database

1. Introduction

Advances in computational methods and technology have facilitated the development of efficient and comprehensive (geo)chemical thermodynamic and physical-chemical models for simulation of the behavior and complex feedbacks of natural systems. Computational thermodynamics has many applications in geochemistry, petrology, chemical engineering, chemistry, and materials research, because multicomponent-multiphase systems can be simulated at pressure-temperature conditions and over timescales that are not accessible to direct observation and laboratory experiments. These simulations are useful for solving environmental problems (e.g. long-term prediction of radioactive waste disposal or contamination of groundwater), designing and improving industrial processes (e.g. formation and stability of different materials), and understanding the evolution of geochemical systems from the surface to the deep Earth. In particular, geochemical reactive-transport simulations that couple thermodynamic fluid-mineral equilibria, kinetics of mineral dissolution and precipitation, and fluid flow in the subsurface have become essential for understanding and predicting the processes in geosystems relevant for carbon dioxide sequestration, exploitation of geothermal energy, and formation of mineral resources (Steeffel and Lasaga, 1994; Steeffel et al., 2005; Xu et al., 2011; Zhang and Parker, 2012; Hoffmann et al., 2012).

One of the main steps in thermodynamic equilibrium modeling is to calculate the molar Gibbs free energy of all components of all phases as a function of temperature, pressure, and composition. The equilibrium state (at fixed composition, temperature and pressure) is then determined by the global minimum of the total Gibbs energy of the system. The Gibbs Energy Minimization (GEM) algorithm (Karpov et al., 1997; Kulik et al., 2013) finds the unknown phase assemblage and speciation of all phases by minimizing the total Gibbs energy of the system while maintaining the mass balance. Conversely, the Law of Mass Action (LMA) algorithm (Reed, 1982) finds the equilibrium speciation by solving a system of nonlinear equations that combine mass balance and mass action expressions. The method directly minimizes the mass balance residuals and performs additional loops if the stable phase assemblage is not known in advance. Although LMA algorithms can perform faster in simple chemical systems, the GEM method is better suited for solving for the equilibrium in complex heterogeneous chemical systems with many non-ideal multicomponent solution phases (Kulik et al., 2013; Leal et al., 2014). The main input and output parameters and properties used in the GEM approach are listed in Table 1.

Table 1. GEM method input and output parameters/properties.

GEM input	GEM output
<ul style="list-style-type: none"> • List of independent components (elements) • List of phases • List of species (dependent components) in all phases • Standard state Gibbs free energy (G°) for each species at T, P of interest • Bulk (elemental) chemical composition • Temperature T and pressure P of interest • Parameters of activity models for components of solution phases 	<ul style="list-style-type: none"> • Chemical system speciation (mole amounts of dependent components in all phases) • Total volume and Gibbs energy of the system • Activity coefficients of dependent components in their respective phases • Amount, volume, mass, and bulk elemental composition of the multicomponent phases • Effective aqueous ionic strength (IS), pH, pe, Eh (in aqueous systems).
<hr/>	
Optional:	
<ul style="list-style-type: none"> • Specific surface areas of phases • Kinetic rate parameters for phases • Additional metastability restrictions for some phase components 	

As any numerical model, a chemical-thermodynamic model is a mathematical formulation that describes the relevant features of a natural system or its parts. Key components of thermodynamic models are the control (input) parameters, which may be empirical (e.g. compositions of fluids and rocks, fluid/rock ratios) or may represent some physical-chemical properties of the system (thermodynamic properties, temperature, pressure). For example, the Helgeson-Kirkham-Flowers equation of state (HKF EoS) (Helgeson et al., 1981; Shock and Helgeson; Tanger and Helgeson, 1988; Shock et al., 1992), which describes the temperature and pressure dependence of the standard-state thermodynamic properties of aqueous species, uses 7 semi-empirical parameters. These parameters are not accessible to direct observation or measurement, but need to be derived by regressing experimental data for measurable quantities (mineral solubility, heat capacity and volume of aqueous solutions).

Regression of model parameters involves adjusting them by iterative numerical methods in such a way that the differences between model-calculated properties and their experimentally determined

counterparts become minimal. Input parameters derived using different models and separate sets of the experimental data may not be mutually consistent and will lead to unrealistic model predictions. Simultaneous processing of large sets of experimental data and regressing them with the same thermodynamic model ensures that the derived thermodynamic properties are internally consistent and accurately reproduce the experimental data (Anderson and Crerar, 1993). To make this possible, a single code framework must integrate a comprehensive collection of appropriate thermodynamic models, an efficient numerical method to perform the thermodynamic computations, a numerically stable and efficient optimization algorithm to adjust the parameter values, and an extensible, flexible collection of the experimental data and parameter optimization tasks.

GEMSFITS is a code package that can adjust any input parameters or properties (see Table 1) for GEM-based modeling of (geo)chemical equilibria, provided that experimental datasets are available and that the input parameters are sensitive to them. This is a much more extended scope than that of the previous prototype GEMSFIT (Hingerl et al., 2014), which was aimed only at fitting interaction parameters of thermodynamic activity models. The GEMSFITS codes are fully compatible with the GEM-Selektor software package (<http://gems.web.psi.ch/GEMS3>) and the GEMS3K numerical kernel (Kulik et al., 2013). The GEMSFITS package consists of `gemsfit2` (parameter optimizer) and `gfshell2` graphical user interface (GUI) codes, which both access the same NoSQL database files. The package offers a general and flexible way for handling the experimental database, setting up parameter optimization or inverse modeling tasks, and visualizing and analyzing the results of the fitting. The `gemsfit2` code can efficiently run complex multi-dimensional parameter optimization problems and can be parallelized. In this paper, we describe the main features and computational methods of the GEMSFITS package and demonstrate the efficiency of its application to chemical-thermodynamic problems of aqueous geochemistry and petrology.

2. Review of software for optimization of thermodynamic model parameters

Generic fitting code packages like MATLAB (MathWorks, 2012), HOPSPACK (Plantenga, 2009), UCODE (Poeter and Hill, 1998), DAKOTA (Eldred et al., 2007), or PEST (Doherty and Hunt, 2010) perform nonlinear parameter estimation using data exchange via input and output files; the user has to manually implement all thermodynamic models involved in the parameter fitting. In many cases, complex scripts are needed to call the objective function subroutine that compares the model output properties with their

empirical (experimental) counterparts. One of the benefits of these methods is that they are independent of the modeling software package.

Another approach is to couple the optimization routine to a chemical solver. The FITEQL code (Herbelin, 1999) uses nonlinear least-square optimization for determining chemical equilibrium constants from experimental solubility and titration data. Due to limitations of the implemented algorithms, the program has convergence issues in more complex fitting problems. Furthermore, because of a lack of normalization options, the fitting results are more sensitive to the data points with high absolute values (Karamalidis and Dzombak, 2010). The program comes together with the MINEQL LMA chemical solver that can be used to perform calculations only in aquatic systems with pure solids, simple ideal solid-solutions, and/or adsorption at low temperature (0-50 °C) and low to moderate ionic strength (<0.5 M) (Westall et al., 1976). This substantially limits the range of possible applications.

The PhreePlot package (Kinniburgh and Cooper, 2011) includes the embedded PHREEQC LMA chemical solver (Parkhurst and Appelo, 2013) and thus has a direct access to all chemical speciation models available in PHREEQC. However, preparing a more complex fitting task in PhreePlot (e.g., when experiments with more than one dependent value are regressed and/or when parameters belong to separate models) is regarded as quite an advanced task (Kinniburgh and Cooper, 2011). The program involves only local non-gradient based (derivative-free) optimization algorithms, and the output summary statistics is minimal (calculation of confidence intervals and sensitivity and parameter correlation analysis are not performed). The possibilities for optimization of several model parameters against large experimental datasets are limited because the code is not parallelized. The LMA algorithm used in PHREEQC also makes the parameter regression for complex non-ideal solution models or non-ideal fluid mixtures not possible.

3. Key features of GEMSFITS

In comparison with the software packages reviewed above, GEMSFITS offers a robust collection of non-linear (global, local, gradient based, and non-gradient based) optimization algorithms, coupled with the efficient GEMS3K chemical solver (Kulik et al., 2013) that includes a broad range of mixing models for solution phases collected in the `TSolMod` library (Wagner et al., 2012). Another relevant feature of GEMSFITS is that it can store, retrieve, and manage in a flexible and efficient way large amounts of structured experimental data in a NoSQL database. This feature also allows the user to produce and store

different optimization tasks using various selections of data sets and experimental samples. In contrast, other available codes (see Section 2) do not provide the convenient tools for managing the experimental data and optimization tasks, which makes it very difficult to develop and maintain comprehensive datasets with experimental data for large chemical systems.

Unlike the earlier prototype GEMSFIT (Hingerl et al., 2014) that used an structured query language (SQL) database server, the GEMSFITS codes can manage and access the experimental data collected in databases in the industry-winning NoSQL BSON format (Binary JSON JavaScript Object Notation; <http://www.mongodb.com/nosql-explained>; <http://bsonspec.org/>). This difference is important, because in SQL databases, the individual records (e.g., for experiments or samples) are stored similar to rows in tables, whose columns contain properties (e.g., chemical element amount, temperature, pressure) that all must be initially specified. Different data types are stored in separate tables that are connected with the join tables needed to execute complex selection queries (e.g., selecting experiments with a certain phase). The SQL relational database implementation becomes extremely complex when describing hierarchical chemical systems with all properties, phases and components, because it has to rely on dozens of data and join tables.

Conversely, a NoSQL database (as used in GEMSFITS) stores the data as documents, each defined as a JSON object. The simplest object is a { *<key>* : *<value>* } pair, for instance { “temperature” : 298.15 }. The *<value>* can be either a constant (string, number, binary data block), an ordered array of values [*<value1>*, *<value2>*, ...], or a nested key-value pair. This allows describing complex hierarchical data structures in a natural, straightforward and human-readable way, similar to structured types in object-oriented programming languages. Each JSON document can contain any number of nested key-value pairs. Compared to the structure of SQL databases, which require a defined schema before adding the data (i.e., SQL databases need to know all column headers in all tables in advance), the NoSQL databases allow inserting data without a predefined schema because every JSON document (database record) may, in principle, have a different data structure. This approach is much better suited for storing large volumes of weakly structured data without making any changes to the already stored documents, and implementing future extensions of the database. Thus, the NoSQL database can store definitions of complex chemical systems or experimental data, which are flexible and easy to process in codes based on object oriented programming.

Based on our experience with the earlier prototype GEMSFIT (Hingerl et al., 2014), the GEMSFITS package has been completely redesigned to efficiently solve the following classes of parameter optimization problems (combinations are possible as well):

- 1) Fitting of interaction parameters of mixing models including aqueous activity models;
- 2) Optimization of thermodynamic properties such as standard state Gibbs free energy $\Delta_f G_{298}^0$ of compounds, or equilibrium constants of chemical reactions;
- 3) Thermobarometry (finding temperature and pressure of formation for the known phase speciation);
- 4) Inverse titrations (e.g. finding the bulk composition that results in prescribed pH);
- 5) Combined (nested) titration-solubility, titration-adsorption and similar fitting problems.

Thus, the GEMSFITS package has the capability to optimize any GEMS3K input properties or parameters (see Table 1). Several parameters from different groups can be fitted simultaneously, either as unconstrained parameters, or some of them can be bounded, reaction-constrained or linearly-constrained. GEMSFITS can perform extensive statistical evaluation of the fitted parameter values, thus helping the user to evaluate the quality and physical significance of the regression results. The program can also be executed on parallel computer architectures, reducing the amount of computing time, and making it possible to run very large optimization problems with the Monte Carlo generated statistics. The GEMSFITS package is available for free download from <http://gems.web.psi.ch/GEMSFITS>, eventually open-source.

4. Methods and data

The GEMSFITS package is composed of the `gemsfit2` parameter optimizer and the `gfshell2` graphical user interface (GUI), both accessing the same NoSQL database in JSON/BSON (EJDB, <http://ejdb.org>) format where the experimental data, the fitting task definitions, and the task calculation results are stored. This functionality enables the user to access and manage the database, to specify the fitting tasks by editing their definitions, to run the parameter optimization (optionally generating statistics of the fitted model parameters), to view, and to plot and print fitting results and statistics. In this way, an efficient and flexible workflow of GEM input parameter optimization is made possible.

The setup of the fitting task (to be performed by the `gemsfit2` code) is provided in the task input specification file that can be exported from the `gfshell2` or prepared using any external text editor. The

task setup controls the selection of experimental data from the database, defines what model output values are retrieved, and how they will be compared within the global and/or nested objective function. Furthermore, it selects the model input parameters that should be fitted (and provides initial values, optionally lower and upper bounds), and defines the choice of the optimization method, the sample weighting rules, and the options for statistics. Upon execution, the `gemsfit2` code reads the task specification input file, the GEMS3K chemical system definition files, and the NoSQL database with experimental data, performs the requested calculations (writes the steps into a log file), and finally writes the results into output files in comma-separated values (csv) format (results can be imported into NoSQL database in connection with the fitting task that generated them).

4.1. Architecture of GEMSFITS

The GEMSFITS software package consists of four main components:

- 1) The NoSQL database (collection of documents describing experimental samples, fitting task definitions, and fitting task results) in BSON/JSON format, with tools for the database management;
- 2) The `gemsfit2` parameter optimization code that reads the experimental data from the database, executes a “fitting task” as described in the task specification file (exported from the task definition database record, or edited separately), and writes results into a set of csv format output files;
- 3) The GEMS3K chemical solver (including the `TSolMod` library of mixing models) (Wagner et al., 2012; Kulik et al., 2013) embedded into the `gemsfit2` code, which reads the chemical system definition from the set of GEMS3K input files provided, and calculates chemical equilibria and/or phase activity models whenever called by the parameter optimizer;
- 4) The `gfshell2` GUI with graphics widget and help viewer that assist the user in accessing the database, preparing the input files, running the fitting tasks, and exploring the results. It keeps track of the task definition (which generates the task specification file) and task results records in the database, and provides text editors and graphics for input and result data.

The GEMSFITS GUI-based workflow is organized in projects consisting of one or more fitting tasks, as illustrated in Fig. 1. A project defines one general parameter optimization application. Each project refers to an experimental database file and a set of GEMS3K chemical system definition files. A fitting task is defined by a task specification record containing all the settings for the optimization process. The user can

configure alternative task definitions with different optimization options, save task definitions to the database, run these different optimization tasks, and view and save their results.

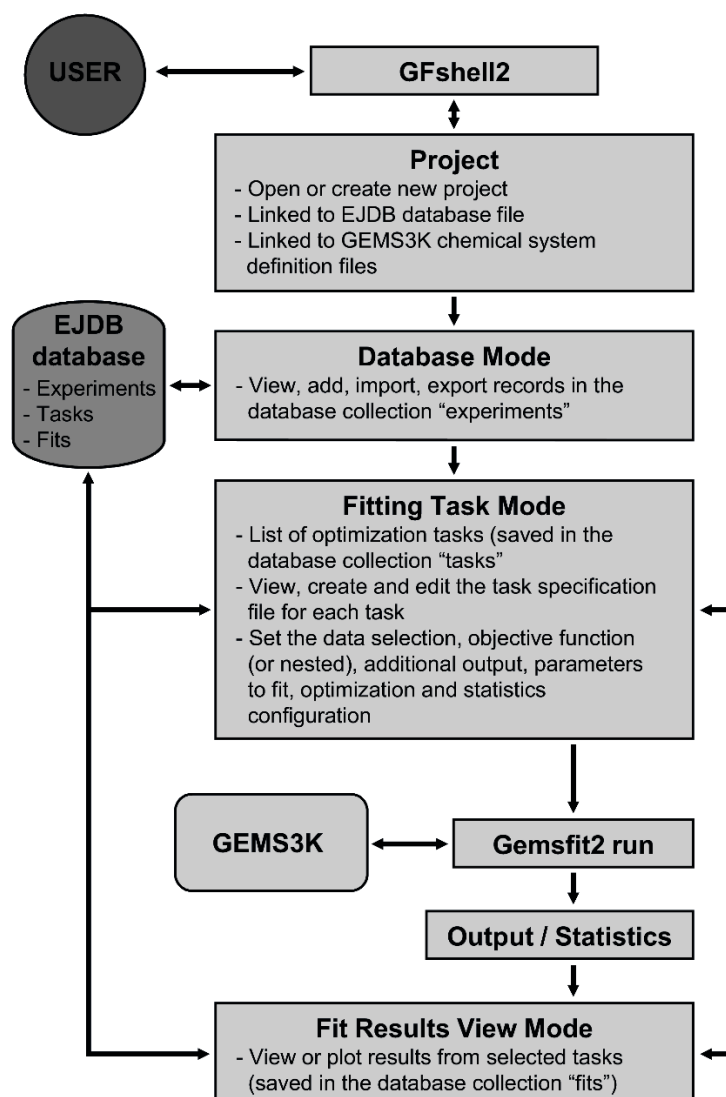


Fig.1. Flow chart of the `gfshell2` code illustrating the four main modes of code operation (Project Mode, Database Mode, Fitting Task Mode and Fit Results View Mode) and the links to the EJDB database.

4.2. Experimental and task database management

Currently, the NoSQL database is implemented as files within the local computer application using the EJDB library (Embedded JSON Data Base engine, <http://ejdb.org>), which is a lightweight variant of the MongoDB database (<http://mongodb.org>). GEMSFITS databases can be transferred to a MongoDB server in the future, if required for distributed database development projects. The database is initially created, extended, or updated by importing the experimental data in csv format, using the `gfshell2` menus. The csv files can be exported from pre-formatted spreadsheets.

The database operations can be performed directly in the database using the `EJDB console` (<http://ejdb.org>) or via the `gfshell12` GUI. The latter allows also for exporting the data into JSON format files. The database stores the data in a BSON (Binary JSON) format (<http://bsonspec.org/>).

An advantage of the JSON format is that the data is represented through a hierarchical structure as { <key>: <value> } pairs, such as { "phase": "aq_gen" }. In place of <value>, an array of ordered values [<value1>, <value2>, ...] can be used, e.g. { "phase": ["aqueous", "gaseous", "calcite"] }, or a subordinated { <key>: <value> } pair, for instance { "phase": { "aggrstate": "electrolyte", "name": "aqueous" } }. This recursive notation can represent any hierarchical, structured, and ordered data objects such as those used in advanced object-oriented programming languages. This format makes it easy to handle data that describe various experimental settings and chemical systems, allowing for streamlined expansion of new database documents (records) without modifying the already existing data records or the entire database architecture.

The `EJDB` stores all data records ('documents') in 'collections'. A collection is a group of documents having similar (but not necessarily the same) structure. In a `GEMSFITS` project database, the experimental data is stored in the collection "experiments". The optimization task specification is saved in the "tasks" collection; after each successful optimization run, the results can be saved in the "fits" collection of the project database file. In this way, the data records can be viewed, edited and used at any later time. Experimental data records are identified by two keywords "expdataset" and "sample", a combination of which should be unique in the entire database collection, and forms a 'data record key'. The "expdataset" value (string) typically refers to a paper or a report, or part of it describing a substantially different setup of the experiments. The "sample" value (string) refers to a single sample or a measurement point with defined temperature, pressure, and composition. The experimental data consists usually of three subsets or sections, which are 'Experiment description', 'System composition', and 'Measured results' (Table 2).

The major advantage of using the `EJDB` with standalone files is that these files are part of the `GEMSFITS` application and can be easily packaged together with the rest of the fitting project, without requiring access to a database server. Because `EJDB` uses the same BSON C Application Programming Interface (API) as the `MongoDB`, the data can be backed-up to JSON format files (<http://www.json.org/>) and restored from them into a `MongoDB` server (Plugge et al., 2010), if necessary.

Table 2. Main types of experimental data that can be added to the database and used in calculation of the objective function for the sum of residuals (Measured results).

Experiments description	System composition	Measured results
<ul style="list-style-type: none"> • Experimental dataset (expdataset) • Sample name (sample) • Comment • Temperature (sT) • Pressure (sP) • Volume (sV) of the system 	<ul style="list-style-type: none"> • Chemical composition of the system in formula units ($comp$) • Upper and lower additional metastability constraints for phases or species (UMC, LMC) • Unit of measurement and estimated error for all entries 	<ul style="list-style-type: none"> • Concentration of elements (independent components IC) in aqueous phase • Mole amount of independent components (IC) in phases-solutions • Mole ratios (MR) of elements in gaseous, solid, melt, and aqueous phases • Phase properties: mass (Q), volume (pV), excess Gibbs energy of mixing (Gex), density (RHO), pH, pe, eH, ionic strength (IS), alkalinity (alk), surface area ($sArea$), osmotic coefficient of water ($oscw$) • Amounts or concentrations of dependent components (DC) in non-aqueous phases • Units of measurement and estimated experimental errors for all entries

4.3. Parameter optimization code

The `gemsfit2` code performs the actual parameter optimization as described in the task specification file, according to the flow chart shown in Fig. 2. The program reads in the experimental data from the database according to the "DataSelect" query. Using the GEMS3K input chemical system definition files, `gemsfit2` creates an array of chemical system definitions in the computer memory, one per 'experimental sample'. At the next step, the program checks if there is any nested objective function defined in the "DataTarget" section of the task specification file. If no nested function is defined, the program proceeds with computing the equilibrium state for each given sample by calling the embedded GEMS3K chemical speciation solver. For certain objective functions, an alternative to computing the equilibrium is to use direct access to the `TSolMod` code library (Wagner et al., 2012) of geochemical-thermodynamic mixing models. This significantly reduces computational time in cases when the equilibrium phase composition is known and the calculation of phase equilibrium is not necessary (e.g., optimizing activity model parameters against data such as osmotic coefficients of aqueous solutions or excess Gibbs energies of a solid-solution).

The sum of residuals is computed using the objective function terms that are defined in the "DataTarget" section. Each term describes what measured data should be compared with the computed counterpart (e.g., measured concentrations of dissolved elements like Al and Si in solubility experiments).

The sum of residuals can be computed as a classical “sum of squares” or other implemented variants, and several alternative weighting methods can be applied (Table 3). The sum of residuals (the “target”) is then transferred to the chosen optimization algorithm, which will generate new input parameter values, trying to minimize the sum of residuals. This loop consists of computing the equilibria (or phase property) for all samples, calculating the sum of residuals, and refining the fitting parameters. It is repeated until the defined threshold for convergence is reached. At the end of the optimization process, the output (results and statistics) is written into csv formatted text output files.

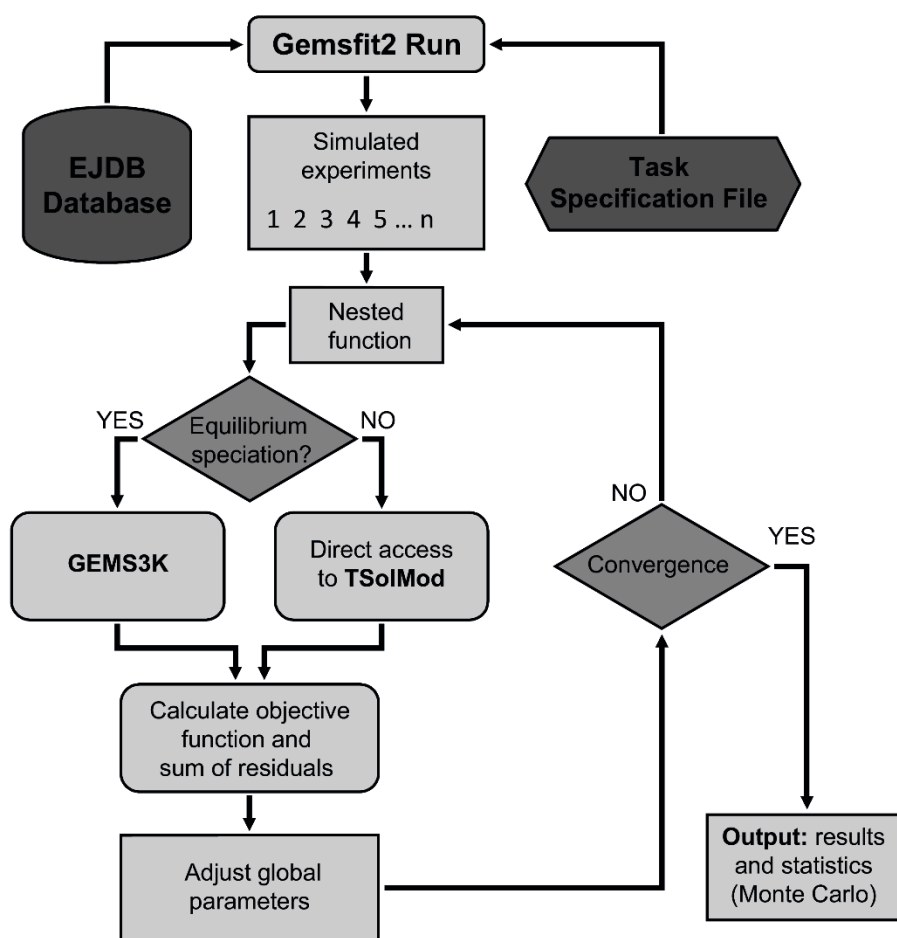


Fig. 2. Flow chart of the `gemsfit2` global optimization loop illustrating the connections to GEMS3K and TsoIMod code libraries.

The nested objective function, for the first time implemented in `gemsfit2`, allows the possibility to set up a fitting task in which the GEM input properties for each experiment (such as temperature, pressure and bulk elemental composition) must always be adjusted against measured data such as pH, alkalinity, or fugacities of gases (Fig. 3). Typical examples of such experimental data include pH-dependent solubility

data for minerals, or pH edges for ion adsorption, where pH is measured, but not explicitly defined in the experimental sample compositions. Most GEM algorithms do not allow to set pH, alkalinity or fO_2 as a direct input. For a given system definition, these activity-based measurable output properties can only be adjusted to desired (measured) values by changing the bulk composition (e.g., titration with acid or base), temperature or pressure.

Table 3. Functions for calculating the residuals and weights in `gemsfit2`.

Residual functions and weights
$F = \sum_{i=1}^n w_i (f_i - y_i)^2 w_o w_e w_{Tu}$
$F = \sum_{i=1}^n w_i \left(\frac{f_i}{\bar{y}} - \frac{y_i}{\bar{y}} \right)^2 w_o w_e w_{Tu}$
$w_i = \frac{1}{\sigma_i}$
$w_i = \frac{1}{\sigma_i^2}$
$w_i = \frac{1}{y_i^2}$
$w_i = \frac{\bar{y}}{\sigma_i^2}$
$w_i, w_o, w_e, w_{Tu} = 1 \text{ (default)}$

Properties: f : computed property; y : measured property; \bar{y} : measured property average; σ : error; n : number of experiments; F : sum of residuals; w_i : weight related to error or measured property value; w_o : weight related to an objective function term (e.g. measurements of Si concentration can have a different weight than the ones of Ca concentration); w_e : individual experiment weight.

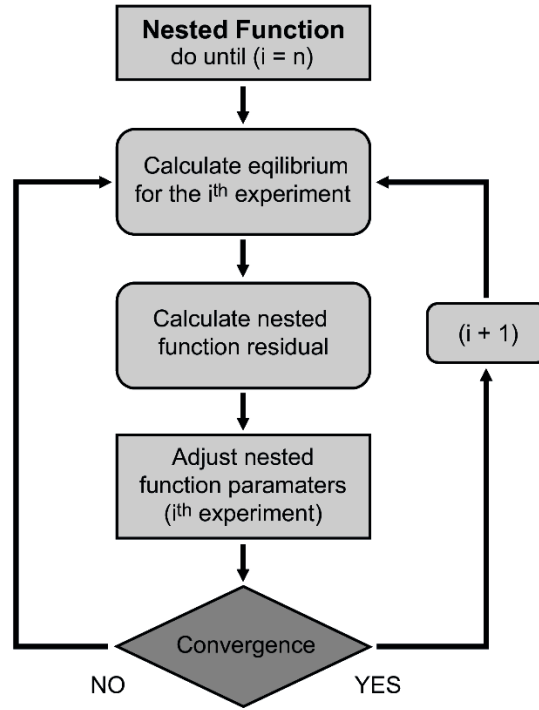


Fig. 3. Flow chart of the `gemsfit2` nested objective function loop (inverse modeling).

If a nested objective function is defined in the "DataTarget" section of the task specification file, the program first loops through all simulated experiments (samples). For each experiment, the program calculates the absolute difference between the measured and the computed property specified in the nested function (any GEM output with experimental counterpart, e.g., pH). The optimization algorithm then adjusts the parameters involved in the nested function until the computed properties achieve the best agreement with the measured ones. When the optimization algorithm converges, the same procedure is done for the next simulated experiment, until the end of the list of experiments. After executing the nested function optimization for all simulated experimental samples, the program proceeds with the main loop for optimizing the global fitting parameters using the top-level objective function (Fig. 2). If no top-level objective function is defined, the program writes the results of the nested function optimization into a csv file as the inverse modeling result.

Parameters that can be optimized by the `gemsfit2` code are summarized in Table 4. Parameters that influence the computed properties in many experiments are optimized using the global objective function (with one or more terms) by minimizing the total sum of residuals

$$\min (total S_{res}) = \sum_0^M F_0 \quad (1)$$

where F_0 is the sum of residuals calculated as described in Table 3, and M is the number of terms in the global objective function.

Table 4. Main classes of adjustable parameters for fitting tasks and parameter fitting modes.

Parameter	Fitting mode	Objective function
• Standard state Gibbs free energy (25 °C and 1 bar) of phase components “ $G0$ ”	• Free (F) • Reaction constrained (R) • Set (S)	Global
• Mixing models interaction parameter coefficients “ PMc ”, “ DMc ” • DQF parameters of end members pure gas fugacities “ $JDQF$ ”	• Free (F) • Set (S)	Global
• Temperature and pressure “ TK ”, “ P ”	• Free (F) • Set (S)	Nested
• Element bulk composition “ bIC ”	• Independent (F) • Linearly constrained (L) • Set (S)	Nested

Parameters specific to an individual experiment only (inverse modeling) are optimized using the nested objective function, by minimizing the absolute residual values of the difference between the computed property f and the measured property y

$$\min(res_i) = \min|f_i - y_i| \quad (2)$$

Three different parameter-fitting modes can be used, which are free (independent) (F), reaction constrained (R), and linearly constrained (L) fitting. Parameters marked as free ‘F’ are optimized independently of each other. Some standard-state molar Gibbs energies of components (25 °C and 1 bar), marked as ‘R’, can be optimized using an additional reaction constraint. This means that at each optimization step, the new value is re-calculated using a (user-provided) reaction equilibrium constant K and the independently optimized values of molar Gibbs energies of any other species that take part in the reaction:

$$Rp = -RT \ln K + \sum_{j=1}^{Nr} P_j c_j \quad (3)$$

Here, Rp is the value of the reaction-constrained parameter, R is the gas constant, T is the temperature in Kelvin, K is the equilibrium constant, Nr is the number of parameters involved in the reaction other than the constrained one, P_j is the value of the parameter involved in the reaction (can be reaction constrained, freely fitted, or fixed), and c_j is the reaction stoichiometry coefficient. The reaction always involves one species constrained by the equilibrium constant, and several others that are independently fitted or have been fixed.

To give an example, the solubility data for halite (NaCl) could be fitted by regressing the properties of aqueous Na^+ and Cl^- species, but adjusting the properties of the aqueous species NaCl^0 through equilibrium constants obtained for the ion association reaction ($\text{Na}^+ + \text{Cl}^- = \text{NaCl}^0$) from independent sets of experiments (e.g., conductance data). Note that temperature and pressure cannot be adjusted simultaneously with implicitly T and P dependent parameters such as "G0" or "PMc". Bulk composition parameters, if marked with 'L', can be linearly constrained to independently fitted composition variables, in order to reproduce the stoichiometry of compounds such as titrants.

An additional "Set" mode (S) is available, to set the input parameter to a new value (different from that initially given in GEMS3K input files, perhaps already fitted in a previous step) to be kept constant during the optimization procedure. This limits the applicability of that new parameter value to a given task while not affecting other tasks.

4.4. GEMS3K chemical solver and its input data

GEMS3K (Kulik et al., 2013) is a chemical equilibrium speciation solver that performs modeling of multiphase-multicomponent geochemical equilibria using a GEM algorithm. GEMS3K provides a built-in selection of equation-of-state and activity models for multicomponent phases available in the `TSolMod` class library (Wagner et al., 2012). The chemical system definition (CSD) and the thermodynamic data arrays are usually set up and exported to GEMS3K input files using the GEM-Selektor v.3 (GEMS3) code package. When setting up the chemical modeling project in GEM-Selektor GUI, the user is guided through a stepwise wizard where the elements comprising the system compositions, the thermodynamic database, and the thermodynamic models of mixing in aqueous and gaseous phases are selected. When defining the chemical system (including components, species and phases), the user can select for each phase one of several non-ideal mixing and activity models or equations of state, depending on the thermodynamic framework of the model. The `TSolMod` class library contains most of the commonly used mixing models for aqueous solutions, gas mixtures, fluid mixtures and solid-solutions (Wagner et al., 2012).

Internally, the `gemsfit2` code calls GEMS3K in the '-init' mode to read GEMS3K CSD input files for creating an array of nodes for chemical systems representing the selected experimental samples, to further simulate these experiments. During the parameter optimization, the chemical equilibrium state and speciation are calculated for each sample using the GEMS3K solver, or the solution models are directly

accessed and the phase properties calculated without computing the complete equilibrium. This is all done based on the chemical system definition specified in the GEMS3K input files.

4.5. GUI and help functionality

The graphical user interface `gfshell2` is designed to assist the user throughout the complete workflow of parameter optimization, including editing of the database, setup of the fitting task, running the regression procedure, and analyzing the results. The `gfshell2` operates in ‘Database Mode’, ‘Fitting Task Mode’, and ‘Fit Results View Mode’ (Fig. 1). Any fitting project created in `gfshell2` GUI refers to an experimental database file and to a set of GEMS3K chemical system definition files.

In *Database Mode*, the experimental records from the project database can be viewed in a list and edited in JSON format (Fig. 4). Bulk data can be backed up and restored from JSON format files or imported and exported from preformatted spreadsheet files previously saved in csv format. Database selection queries can be added and edited in the “Search Query Editor”. The record (sample) list retrieved from the database by a search query can be then exported into the task specification file with the help of a menu command.

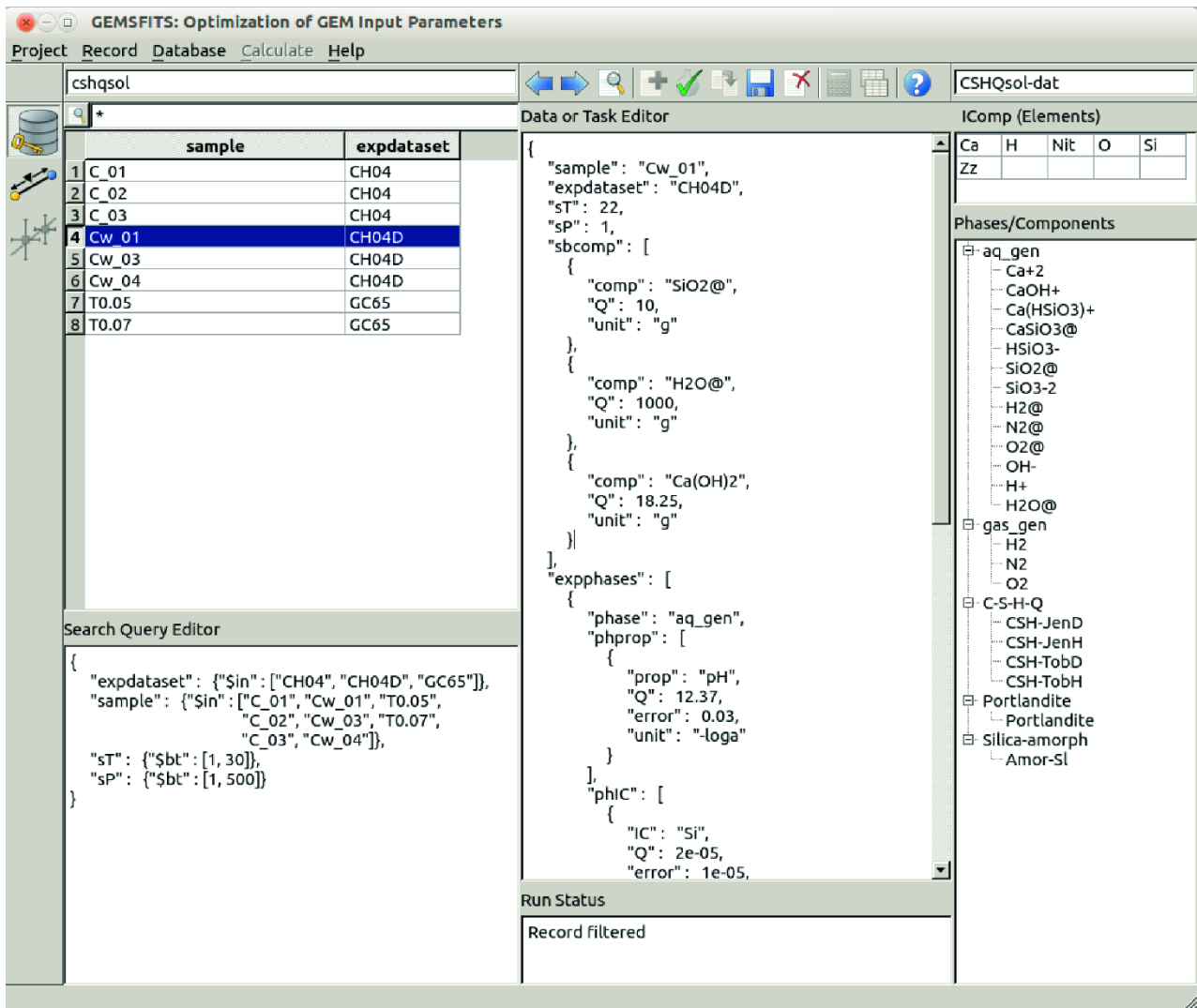


Fig. 4. Screen image of the `gfsHELL2` graphical user interface illustrating the appearance of the Database Mode. The view shows the query result that was retrieved from the experimental database and one sample entry in JSON format. The column on the right-hand side shows the chemical system definition (lists of components, phases and species) that is read from the GEMS3K files.

In *Fitting Task Mode*, the task specification file is viewed and edited (Fig. 5). The specification file contains all the options related to one fitting run. The task specification is viewed in JSON format, where each key has a value and the options. For instance, the pair `"DataSelect": "search_query"` has the value represented by the database search query, or the key-value pair `"OptAlgo": "GN_ISRES"` defines the type of optimization algorithm to be used. Fitting task specifications can be edited, saved to database, or new tasks created from them, and they can all be exported to task input files and executed in the `gemsfit2` code.

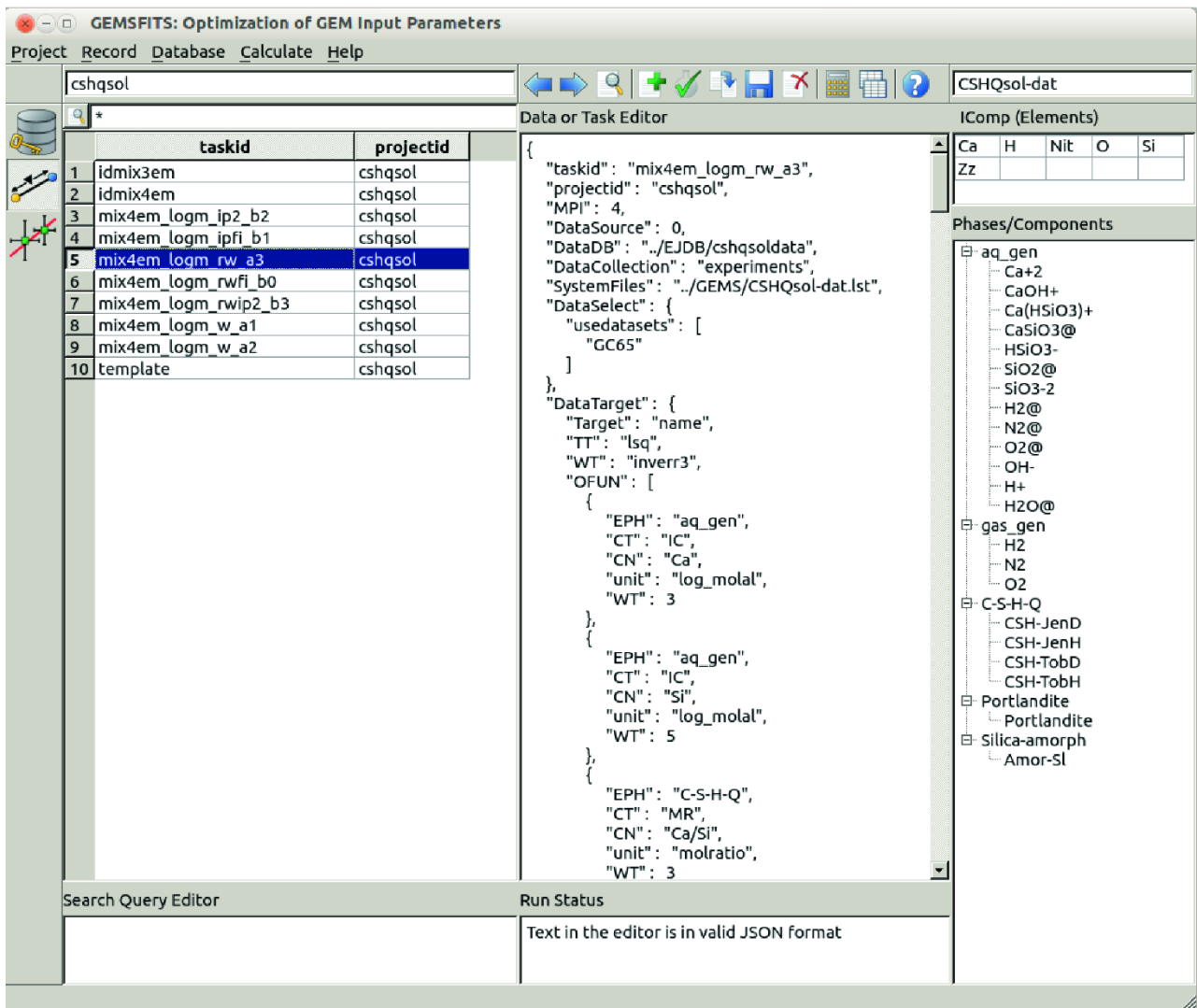


Fig. 5. Screen image of the `gfsHELL2` graphical user interface illustrating the appearance of the Fitting Task Mode. The view shows the list of fitting tasks, the specification file for one task in JSON format and the chemical system definition.

In the *Fit Results View Mode*, the results can be imported from `gemsfit2` output csv files, viewed, and plotted. This functionality is useful for rapid evaluation of the quality of the fitting results. In this mode, the `gemsfit2` output can be displayed in several tabs in a spreadsheet format (Fig. 6). The output contains:

- 1) *Fitted Parameters* with the initial and fitted values of the parameters and the associated parameter statistics;
- 2) *Fit for Samples* containing the dependent/calculated and measured properties, residuals, as well as the weights and other user pre-defined output (Fig. 6);
- 3) *Sum Statistics* with the summary statistics describing the quality of the fit;

- 4) *Sensitivity Data* containing the calculated sensitivities for each measured value to each regressed parameter;
- 5) *Quantile Residuals* with the ordered residuals and their respective quantiles (points taken at regular intervals form the cumulative distribution function of the residuals);
- 6) *MC results* containing the resulting parameter values for each Monte Carlo generated optimization run;
- 7) *Inverse Modeling Results*;
- 8) *gemsfit2 log* containing some diagnostic information related to the fitting run.

sample	expdataset	sT	sP	aq_gen.Ca.meas(y)	aq_gen.Ca.calc(y)	residual	aq_gen.Si.meas(y)	aq
0	T0.05	GC65	25	1	-3.1307683	-3.0186821	-0.11208613	-2.636388
1	T0.07	GC65	25	1	-3.0315171	-3.018685	-0.012832013	-2.5850267
2	T0.14	GC65	25	1	-2.8326827	-3.0186799	0.18599728	-2.4685211
3	T0.23	GC65	25	1	-2.7746907	-3.0186781	0.24398738	-2.3635121
4	T0.32	GC65	25	1	-2.7212464	-3.0186777	0.2974313	-2.3904056
5	T0.41	GC65	25	1	-2.7099654	-3.0186786	0.30871317	-2.4089354
6	T0.43	GC65	25	1	-2.7055338	-3.0186784	0.31314459	-2.412289
7	T0.68	GC65	25	1	-2.69897	-3.050869	0.35189903	-2.4236586
8	T0.76	GC65	25	1	-2.756962	-3.1141873	0.35722531	-2.8326827
9	T0.93	GC65	25	1	-2.6716204	-2.8781244	0.20650399	-3.60206

Fig. 6: Screen image of the `gfshell12` illustrating the appearance of the Fit Result Viewing mode with different spreadsheet tabs (Fitted Parameters, Fit for Samples, Summary Statistics, Sensitivity Data, Quantile Residuals, Monte Carlo Residuals, Inverse Modeling Results, `gemsfit2 log`).

4.6. The `gemsfit2` task specification

The fitting task specification contains all the options and information that are required by `gemsfit2` for performing the computations. It exists as a property-value JSON format editable in `gfshell12`. It has four main sections, which are ‘Data Selection’ (starts with keyword `"DataSelect"`), ‘Data Target’ (with `"DataTarget"`), ‘Parameter Markup’ (with different parameter lists/arrays e.g. `"G0"` for Gibbs free energy parameters; `"bIC"` for elemental bulk composition), and ‘Optimization and Statistics’ (with different options such as `"OptAlgo"` for defining the algorithm; `"StatMCRuns"` for defining the number of Monte Carlo runs). Before running the task, each section must be defined and edited by the user. The query for

selecting the experimental data, the definition of the objective function, and the markup of the parameters to be adjusted are all defined as JSON objects. Using the `gfshell2` program, any task input file can be imported and edited in JSON format, and then saved in the “tasks” collection of the project database. The task specification is automatically exported into the `gemsfit2` input file when the task calculation is started from `gfshell2`.

4.7. Optimization library

Handling complex parameter optimization tasks with numerous parameters and a complex objective function requires a versatile selection of efficient and numerically robust optimization algorithms. Multi-dimensional fitting exercises almost invariably result in convergence difficulties, because of possible local minima and/or highly correlated fitting parameters. For this reason, `gemsfit2` uses the `NLopt` nonlinear optimization library (<http://initio.mit.edu/nlopt>) that was also employed in the earlier prototype (Hingerl et al., 2014). This library provides several global and local minimization algorithms, which can be gradient or non-gradient based (derivative-free). The algorithms implemented in `NLopt` and thus available in `gemsfit2` are listed in Table 5. In the ‘Optimization’ section of the task specification, several options can be used to control the use of `NLopt`, e.g., the type of algorithm, a global upper and lower bound percentage for all parameters, the relative and absolute tolerance for the convergence criterion, the initial step of the parameters, and the maximum number of iterations.

Table 5. Optimization algorithms in `gemsfit2`.

Optimization algorithms (NLopt library)
<p>Global (only bound-constrained problems):</p> <ul style="list-style-type: none"> • GN_ISRES: Improved Stochastic Ranking Evolution Strategy (Runarsson and Xin, 2005) • GN_CRG: Controlled Random Search with local mutation (Kaelo and Ali, 2006) • GN_ESCH: Evolutionary algorithm (da Silva Santos et al., 2010) • GN_ORIG_DIRECT: Dividing Rectangles algorithm (Jones et al., 1993) • GN_ORIG_DIRECT_L: A locally-biased form of the DIRECT algorithm (Gablonsky and Kelley, 2001) • GD_MLSL: Multi-Level Single-Linkage (Rinnooy Kan and Timmer, 1987)
<p>Local:</p> <ul style="list-style-type: none"> • LN_BOBYQA: Bound Optimization By Quadratic Approximation (Powell, 2009) • LN_SBPLX: modified Subplex (Rowan, 1990) algorithm • LN_NEWUOA: using quadratic approximation (Powell, 2004), superseded by BOBYQA (above) • LN_PRAXIS: PRincipal AXIS method (Brent, 1972) • LD_MMA: Method of Moving Asymptotes (Svanberg, 2002) • LD_SLSQP: Sequential Least-Squares Quadratic Programming (Kraft, 1994)

- LD_VAR1: shifted limited-memory variable-metric algorithm (Vlcek and Luksan, 2006)

In the definition of optimization algorithms, G stands for Global, L for Local, N for non-derivative, and D for derivative (gradient-based).

4.8. Weighting and outliers

In many datasets, outlying samples may result in significantly degrading the quality of the fitting and may also cause convergence problems for the fitting algorithms. The simplest way of dealing with extreme outliers is to exclude them from the selection of experiments from the database.

Another option available in `gemsfit2` for moderately outlying samples is to assign the outliers a lower weight. This can be done systematically using a Tuckey’s Biweight function that reduces the influence of outliers (Motulsky and Christopoulos, 2004). The function uses a cutoff value C calculated as the median of all absolute values of residuals multiplied with an arbitrary scaling factor (6 by default) set in the task specification file. Each sample whose absolute value of the residual exceeds C is assigned a weight of 0 (i.e., it will be ignored during the fitting). Other samples are given weights determined by the following equation (where R_i is the absolute value of residual, i is the sample index, and C is the cutoff value as defined before):

$$w_{Tui} = \left(1 - \left(\frac{R_i}{C} \right)^2 \right)^2 \quad (4)$$

Sample weights are recomputed at each iteration of the optimization algorithm. These weights are computed for each different data type that is included in the objective function term (e.g., different medians for dissolved aqueous Al and Si concentrations). Weights can be set for each individual experimental sample, but also for each term of the objective function, as user defined values. More than one weighting option can be selected for one fitting task; Table 3 gives an overview of all possible weighting options.

4.9. Statistics

The majority of the statistics options implemented in `gemsfit2` remain the same (with minor corrections) as in the early prototype GEMSFIT; a detailed summary of them can be found in Hingerl et al. (2014). The main statistics features that can be analyzed are goodness of fit, sensitivity analysis, correlation of parameters (Hill and Tiedeman, 2007), and confidence intervals from Monte Carlo simulations (Motulsky and Christopoulos, 2004). All the statistics options can be set in the ‘Statistics’ section of the task specification.

4.10. Confidence intervals for parameters by Monte Carlo simulation

The philosophy of this method is to generate many datasets (experimental pseudo-data), and perform the optimization of fitting parameters for each generated dataset (Motulsky and Christopoulos, 2004). The resulting distributions of adjusted (fitted) parameters are then used to calculate their standard deviations and confidence intervals. The random scatter is generated as an array of values randomly extracted from a synthetic set of normally distributed data, which have the mean value and the standard deviation equal to that of the residuals resulting from the nonlinear regression. Another option is using a “bootstrap sampling” (DiCiccio and Efron, 1996) of the residuals (useful when the residuals are not normally distributed). This is done by randomly sampling the residuals with the possibility of sampling one residual more than once. The random scatter is then added either to the computed property values or to the respective experimental values, as in eq. (5) or (6) (option set in the Statistics section of the task specification):

$$y_{i,new} = y_{i,old} + s_i \quad (5)$$

$$y_{i,new} = f_{i,old} + s_i \quad (6)$$

Here, y_i represents the measured value, f_i is the computed value, s_i represents the random scatter value, and subscripts *new* and *old* refer to the new ‘synthetic’ measured value and the one used in the actual fitting before the Monte Carlo iterations.

The resulting randomly modified values are then used as new “measured values” for the optimization procedure. The random scatter is computed for each objective function term (different types of data that are included into the objective function) on a normalized scale. For example, one scatter array is computed for all aluminum solubility measurements, and a separate one is computed for the silicon solubility measurements. In the case of bootstrap sampling, the random scatter is computed for each type of the data and the experimental dataset (i.e. one literature reference or set of experiments).

The MC procedure is repeated many times (>100), to make it possible to evaluate the standard deviation of fitted parameter values, which represents the error of the parameters for the given scatter of the experimental data. Symmetric confidence intervals are estimated from standard deviations of parameters by multiplying them with a suitable quantile of Student’s t-distribution (Hill and Tiedeman, 2007).

Due to the large amount of computing required for the Monte-Carlo simulations or, in some cases, for the nested functions and global algorithms, these options may be very time consuming. For this reason, the `gemsfit2` code is parallelized, and it can manage many processing tasks using the OpenMP (Open Multi-Processing) shared-memory multiprocessing library (<http://openmp.org/>). The master process is split into parallel tasks (depending on the type of configuration) that are simultaneously executed, thus dramatically reducing the required computation time.

4.11. Fit-independent statistics

The fit-independent statistics are calculated without invoking the optimized parameter values, using only sensitivities and weights (Hill and Tiedeman, 2007). Sensitivities represent the change in the computed value (y_i) divided by the change in the respective parameter value (b_j). They are calculated using central differences, in which the parameters are both increased and decreased. Mathematically speaking, sensitivity is a partial derivative of a computed value (y_i) with respect to the parameter (b_j). It is approximated from the finite difference:

$$\left(\frac{\partial y_i}{\partial b_j} \right)_b \approx \left(\frac{y_i(b_j + \Delta b_j) - y_i(b_j - \Delta b_j)}{2\Delta b_j} \right) \quad (7)$$

The change in the parameter value is calculated using a perturbator value (Δ), which can be set in the `gemsfit2` task specification file. In most cases, when using a small perturbator value, the calculated sensitivities approach the real sensitivities (Poeter and Hill, 1998), although a too small value could result in insignificant changes in the computed values. It is recommended to investigate the effect of different perturbation values until an optimal one is found. The sensitivities indicate the importance of the observations for determining the estimated parameter.

To better reflect the importance of observations in the parameter optimization, sensitivities have to be scaled. This is because their units are calculated as the measured value divided by the parameter, both of which can have very different units. For this purpose, sensitivities of each observation and parameter are multiplied by the value of the parameter (b_j) and by the weight assigned to the observation (w_i), resulting in the dimensionless-scaled sensitivities (*DSS*):

$$DSS_{ij} = \left(\frac{\partial y_i}{\partial b_j} \right)_b |b_j| w_i^{0.5} \quad (8)$$

Information about the sensitivity of one parameter to the observations is provided by the composite scaled sensitivities (CSS) (Hill and Tiedeman, 2007). They are calculated for each parameter using the dimensionless-scaled sensitivities, as follows (n is the number of observations):

$$CSS_j = \sum_{i=1}^n \left[\frac{(DSS_{ij})^2}{n} \right]^{1/2} \quad (9)$$

The variance-covariance matrix is then calculated as:

$$VarCov(b) = s^2 (X^T w X)^{-1} \quad (10)$$

$$X = \left(\frac{\partial y_i}{\partial b_j} \right)_b \quad (11)$$

where $VarCov(b)$ is a square matrix (the size of which is the number of parameters); s^2 is the error variance, i.e., the weighted squared sum of squared residuals divided by the degrees of freedom (number of observations minus the number of parameters); X is the matrix of sensitivities; and w is the weight matrix.

The correlation coefficient between the j^{th} and the k^{th} parameter is then calculated as follows:

$$PCC(j, k) = \frac{VarCov(j, k)}{VarCov(j, j)^{1/2} VarCov(k, k)^{1/2}} \quad (12)$$

5. Application examples

The practical utility and efficiency of the GEMSFITS code package is demonstrated with examples that represent typical classes of optimization problems in geochemical-thermodynamic modeling. These include (1) fitting of interaction parameters of mixing models, (2) optimization of thermodynamic properties such as the standard molal Gibbs energies of aqueous species or the equilibrium constants of formation reactions, (3) thermobarometry, (4) inverse titrations, and (5) combined nested problems.

5.1. Boehmite solubility and Al speciation

Aluminum is an important element in many rock-forming minerals. Because of its low solubility, Al determines fluid-mineral equilibria and the reaction progress during fluid-rock interaction (Pokrovski, 1998;

Bénézeth et al., 2001; Tagirov and Schott, 2001; Manning, 2006; Mookherjee et al., 2014). Modeling aluminum solubility in geologic fluids has always been problematic and controversial due to the inconsistency of thermodynamic data that involve Al, contradictions between different experimental studies of the pH-dependent solubility of Al minerals, discrepant thermodynamic properties of aluminum-bearing minerals that were used for extracting properties of aqueous Al species from solubility experiments, and because different activity models were used in deriving thermodynamic properties of aluminum aqueous species (Tagirov and Schott, 2001).

In this example, the standard molal Gibbs energies ($\Delta_f G_{298}^0$) of aqueous aluminum species at 25°C and 1 bar were fitted using in situ boehmite solubility experiments performed by Bénézeth et al. (2001). The $\Delta_f G_{298}^0$ of Al^{3+} , AlOH^{2+} , $\text{Al}(\text{OH})_2^+$, $\text{Al}(\text{OH})_3^0$ and $\text{Al}(\text{OH})_4^-$, were simultaneously fitted. The thermodynamic properties of boehmite $\text{AlOOH}_{(\text{cr})}$ were accepted from Verdes et al. (1992) and Hemingway et al. (1991), and these have also been used by Bénézeth et al. (2001) for interpreting their experimental results. Thermodynamic properties of aqueous species at the experimental conditions were calculated in the GEM-Selektor v.3 code using the revised Helgeson-Kirkham-Flowers (HKF) model (Helgeson et al., 1981); water properties were calculated from the Haar-Gallagher-Kell (HGK) model (Kestin et al., 1984). The extended Debye-Hückel aqueous electrolyte model (Helgeson et al., 1981; Oelkers and Helgeson, 1990) was used for calculating the activity coefficients of individual species. The standard state thermodynamic properties and HKF parameters of other species present in the system were taken from Shock and Helgeson (1988) for OH^- , Cl^- , Na^+ , from Tagirov et al. (1997) for HCl^0 , from Sverjensky et al. (1997) for NaCl^0 and from Shock et al. (1997) for NaOH^0 .

The boehmite solubility experiments were performed over a wide range of pH (from 2 to 10) at a salt concentration of 0.03 mol/kg (NaCl), at temperatures between 100 and 290°C, and at saturated water vapor pressures. The experimental method, described in Palmer et al. (2001), has been demonstrated to produce consistent results for numerous samples. The pH of the system was measured in situ using a hydrogen-electrode concentration cell (Bénézeth et al., 2001).

To be able to simulate the exact experimental conditions in the GEM-Selektor v.3 code, the complete chemical composition for each system (experimental sample) must be known. Because concentrations of HCl and NaOH (that were used to adjust the pH in the experiments) were not reported by Bénézeth et al.

(2001), we had to apply the inverse titration approach in GEM-Selektor for adjusting the pH in each experiment (sample) to the experimentally measured value, using the nested objective function as implemented in the `gemsfit2` code. The bulk composition of each simulated experimental sample was first adjusted by adding HCl and NaOH titrants by the free fitting of Na and Cl amounts ('F' type) with linearly constrained O and H amounts ('L' type) to reproduce the titrant stoichiometries.

Using `gemsfit2`, the $\Delta_f G_{298}^0$ values of the aqueous Al species were adjusted and three different fitting cases were considered. In the first case (A), the regression yielded $\Delta_f G_{298}^0$ values for Al^{3+} , AlOH^{2+} , $\text{Al}(\text{OH})_2^+$, $\text{Al}(\text{OH})_3^0$ and $\text{Al}(\text{OH})_4^-$ fitted independently. In the second case (B1), the $\Delta_f G_{298}^0$ of AlOH^{2+} was constrained through the equilibrium constant for the species-forming reaction ($\text{Al}^{3+} + \text{H}_2\text{O} = \text{AlOH}^{2+} + \text{H}^+$). In an additional third case (B2), the initial values for all standard state Gibbs energies were set to 10 kJ/mol more negative, while keeping the same reaction constraint as in the second task (B1). Although the species $\text{NaAl}(\text{OH})_4^0$ is included in the chemical system, it is not controlling solubility at low Na concentrations and the $\Delta_f G_{298}^0$ was therefore fixed at the value from Tagirov and Schott (2001). Initial values of the $\Delta_f G_{298}^0$ and other standard molal thermodynamic properties for the aqueous aluminum species were adopted from Tagirov and Schott (2001). Only the $\Delta_f G_{298}^0$ values were adjusted during regression, while other standard molal properties of aqueous species were kept unchanged.

In all three fitting tasks, the local BOBYQA (Powell, 1994) optimization algorithm was used, and the search bounds for parameters were set to $\pm 5\%$ of the initial value. Initial values, final values, and the associated uncertainties for all fitting runs are given in Table 6. The standard sum of squares function was used with equal weights of 1.0 for all experimental data points:

$$F = \sum_{i=1}^n w_i (f_i - y_i)^2 \quad (13)$$

The solubility data were treated as \log_{10} of the molality. In all runs, the resulting final $\Delta_f G_{298}^0$ values of each species were similar, within their uncertainties (set to two times the standard deviation as estimated using the Monte Carlo method by simulating 1000 random sets of experiments).

In the first task (A), all $\Delta_f G_{298}^0$ values were independently fitted ('F' type). In systems with several independently fitted parameters, where the target function is not highly sensitive to each parameter, typically

some fitted parameters are highly correlated, which makes it difficult to obtain physically meaningful values and to find the optimum solution. In the chemical system investigated, the total solubility is not very sensitive to the independent contributions of Al^{3+} , AlOH^{2+} and $\text{Al}(\text{OH})_2^+$ at pH values below 5, where their predominance fields overlap (Fig. 7). The situation is different at neutral to alkaline pH, where $\text{Al}(\text{OH})_3^0$ and $\text{Al}(\text{OH})_4^-$ dominate the speciation. Uncertainties resulting from the independent fitting show that the $\Delta_f G_{298}^0$ values of all species except AlOH^{2+} could be constrained within reasonable bounds, and that the high uncertainty associated with the $\Delta_f G_{298}^0$ of this species is clearly due to correlations. The resulting correlation coefficients of the independently fitted parameters can be found in Table 7; they show that the correlation coefficient between AlOH^{2+} with Al^{3+} is statistically significant ($r = -0.71$).

Therefore, for the second case (B1), the $\Delta_f G_{298}^0$ value of AlOH^{2+} was reaction-constrained ('R' type) using the equilibrium constants associated with the hydrolysis reaction:



$$\Delta_r G^0 = -RT \ln K \quad (15)$$

Values of equilibrium constants ($\log_{10}K$) for the above reaction for each experimental temperature were taken from Palmer and Wesolowski (1993), who determined them by a potentiometric method. After each optimization iteration, the new value for $\Delta_f G_{298}^0$ of AlOH^{2+} was recalculated using the $\log_{10}K$ of the above reaction and a new independently adjusted $\Delta_f G_{298}^0$ of Al^{3+} . The final results (Table 6) of the two fitting approaches show that the $\Delta_f G_{298}^0$ value of Al^{3+} is quite well constrained.

In the third task (B3), the initial values of $\Delta_f G_{298}^0$ of all Al species were set to 10 kJ/mol more negative than in two previous tasks, in order to investigate the significance of initial guesses for the final regressed parameters; all other settings were kept the same. In the first run, the BOBYQA algorithm converged to a local minimum (Table 6), clearly identified by the sum of squares being two times larger than in the two other tasks. A second run was performed using the final values from the first run as initial values. This time, the algorithm converged to a similar minimum as for the preceding cases, and the final values were almost identical.

The main improvement compared to other datasets for aqueous Al speciation is that the final optimized $\Delta_f G_{298}^0$ values were derived fully consistent with the selected aqueous electrolyte model (extended Debye-Hückel equation), the HKF EoS, the selected thermodynamic properties of boehmite, and all the selected experimental data. This opens up the possibility to derive alternative geochemical-thermodynamic datasets using different standard state data for the solubility-controlling mineral phases, different sets of aqueous species, different activity coefficient models (e.g., derive standard state properties based on a Pitzer model), or different experimental datasets. Furthermore, if new high-quality experimental data would become available in the future, thermodynamic properties of the aqueous Al species could be re-derived using the same approach. The key message here is that GEMSFITS optimized values will always be consistent with the selected models and input data, best reproducing the experimental data. Figure 7 shows \log_{10} activity values of aluminum species and total aluminum concentration as function of increasing pH, comparing directly the values calculated with the initial thermodynamic data and with the final adjusted ones. Experimental data points are plotted for comparison, and final regressed $\Delta_f G_{298}^0$ values obviously show a much better agreement with the experimental Al solubility data.

Table 6. Initial and regressed (final) values of $\Delta_f G_{298}^0$ (kJ mol⁻¹) for selected aluminum species. The uncertainty (2σ) represents 2 times the standard deviation estimated from 1000 Monte Carlo runs.

Species	(1)Initial $\Delta_f G_{298}^0$	(2)Final $\Delta_f G_{298}^0$	(2)Error (2σ)	(3)Final $\Delta_f G_{298}^0$	(3)Error (2σ)	(4)Final $\Delta_f G_{298}^0$	(5)Final $\Delta_f G_{298}^0$
AlOOH(cr)	-917.82	-917.82	1.9	-917.82	1.9	-917.82	-917.82
Al ³⁺	-487.5	-486.8	1.4	-486.3	0.5	-486.4	-486.4
AlOH ²⁺	-696.3	-695.8	10.0	-695.1	0.7	-695.2	-695.3
Al(OH) ₂ ⁺	-899.5	-897.9	1.2	-898.6	1.0	-897.1	-898.0
Al(OH) ₃ ⁰	-1101.7	-1105.0	1.0	-1104.6	0.9	-1106.2	-1104.9
Al(OH) ₄ ⁻	-1305.8	-1307.2	0.4	-1307.2	0.4	-1303.9	-1307.2
NaAl(OH) ₄ ⁰	-1567.4	-1567.4		-1567.4		-1567.4	-1567.4
Sum of squares	23.03	13.99		13.97		27.12	13.96

(1)The initial $\Delta_f G_{298}^0$ value for AlOOH(cr) was adopted from Verdes et al. (1992) initial values for all aqueous Al species were taken from Tagirov and Schott (2001). The $\Delta_f G_{298}^0$ values for AlOOH(cr) and NaAl(OH)₄⁰ were fixed in the regression.

(2)Case A, the $\Delta_f G_{298}^0$ of all species was fitted independently.

(3)Case B1, the $\Delta_f G_{298}^0$ of Al³⁺ was fitted directly, while the $\Delta_f G_{298}^0$ of the AlOH²⁺ species was constrained by applying the equilibrium constant for the species-forming hydrolysis reaction (Al³⁺ + H₂O = AlOH²⁺ + H⁺). For this species, the

2σ error was calculated from the error of the $\Delta_f G_{298}^0$ for Al^{3+} and the error of the equilibrium constant for the species-forming reaction.

⁽⁴⁾Case B2, fitted by setting all initial values of $\Delta_f G_{298}^0$ smaller with 10.0 kJ/mol. The $\Delta_f G_{298}^0$ of Al^{3+} was fitted directly, while the $\Delta_f G_{298}^0$ of the AlOH^{2+} species was constrained by applying the equilibrium constant for the species-forming reaction.

⁽⁵⁾Case B2, same as ⁽⁴⁾ and using as starting values the final values of ⁽⁴⁾.

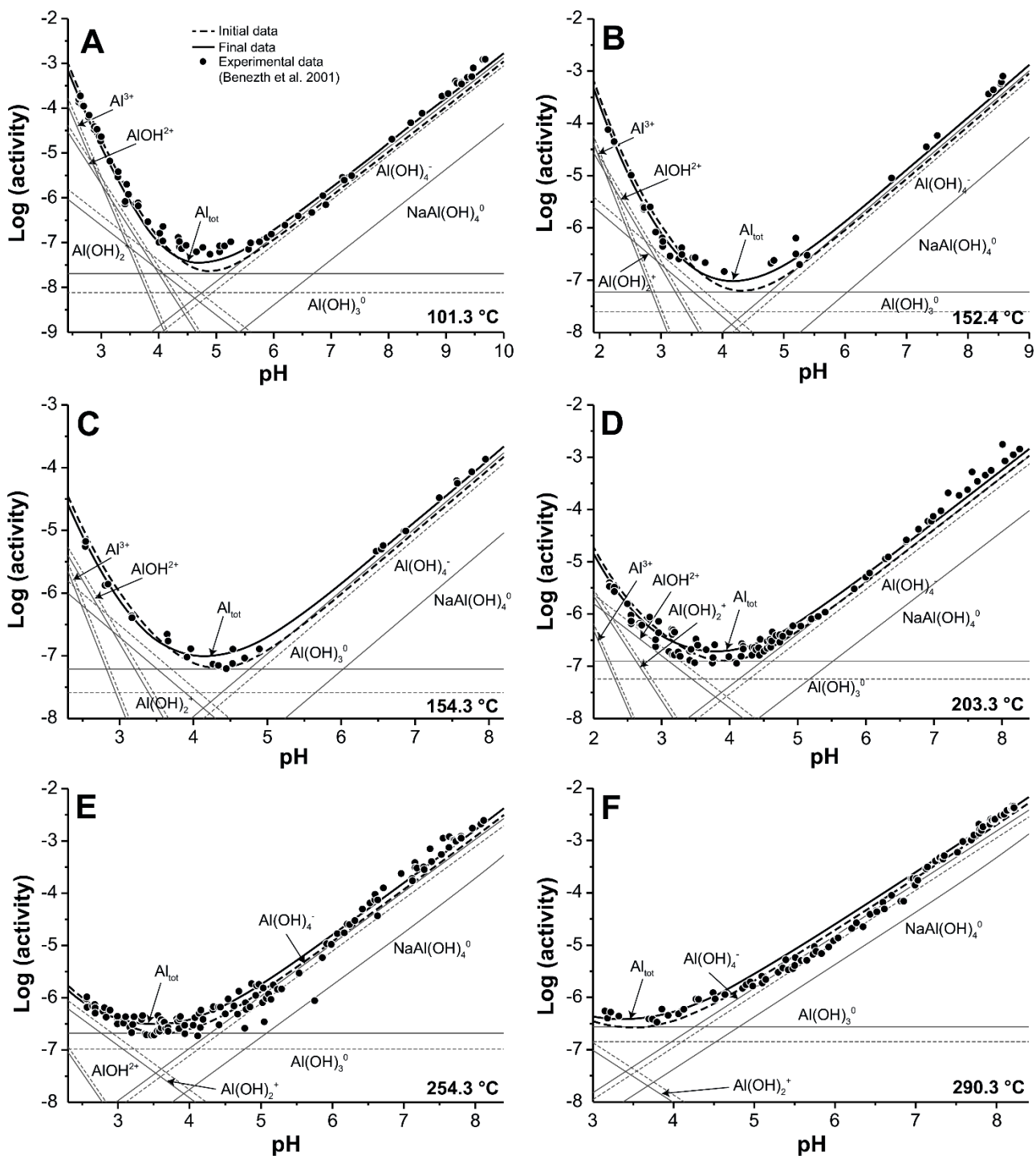


Fig. 7 Plots showing the aluminum solubility and speciation as function of increasing pH and the total aluminum concentration, calculated with GEM-Selektor before (dashed lines) and after the optimization (solid lines) of the standard state properties of the aqueous Al species Al^{3+} , AlOH^{2+} , $\text{Al}(\text{OH})_2^+$, $\text{Al}(\text{OH})_3^0$ and $\text{Al}(\text{OH})_4^-$. Filled circles represent the experimental data points. Plots are shown for temperatures of (A) 101.3 °C, (B) 152.4 °C, (C) 154.3 °C, (D) 203.3 °C, (E) 254.3 °C and (F) 290.3 °C.

Table 7. Parameter correlation matrix for the final regressed $\Delta_f G_{298}^0$ values of Al species.

	Al ³⁺	AlOH ²⁺	Al(OH) ₂ ⁺	Al(OH) ₃ ⁰	Al(OH) ₄ ⁻
Al ³⁺	1.0	-0.71	0.30	-0.11	-0.02
AlOH ²⁺		1.0	-0.56	0.25	0.04
Al(OH) ₂ ⁺			1.0	-0.65	0.16
Al(OH) ₃ ⁰				1.0	-0.29
Al(OH) ₄ ⁻					1.0

The correlation matrix was calculated for the case where $\Delta_f G_{298}^0$ values of all species were fitted independently.

5.2. Ti in quartz: Solid-solution geothermometry

Several authors studied the temperature and pressure dependence of the substitution of Ti for Si in quartz (Ostapenko et al. 1987; Wark and Watson, 2008; Thomas et al., 2010; Huang and Audétat, 2012). This solid-solution is commonly applied as a thermobarometer (Rusk et al., 2008; Smith et al., 2010; Wilson et al., 2012; Kidder et al., 2013).

The objective of this example was to test the optimization of mixing model parameters for solid-solutions using GEMSFITS. Therefore, a regular SiO₂-TiO₂ (quartz-rutile) solid-solution model was constructed in the GEM-Selektor database. Then three coefficients (constant term and linear terms for temperature and pressure dependence) of the regular interaction parameter were fitted against the experimental data (Wark and Watson, 2008; Thomas et al., 2010). The optimized solid-solution model was then used in a GEMSFITS inverse modeling task aimed at determining the temperature of quartz crystallization using measured Ti concentration data in natural quartz samples from Kidder et al. (2013) along with those from low-grade metamorphic quartz veins in accretionary-wedge sediments of the Swiss Alps (Miron et al., 2013).

The solubility of Ti in quartz in the presence of pure rutile was experimentally measured between 600 and 1000 °C and between 5 and 20 kbar using a piston-cylinder apparatus (Wark and Watson, 2008; Thomas et al., 2010). A total number 31 data from these experimental studies were added to the GEMSFITS TiQ database. The chemical system was set up using GEM-Selektor and exported to a set of GEMS3K input files. The Peng-Robinson multicomponent fluid model (Anderson and Crerar, 1993) was used to describe the fluid phase (with fluid species H₂O, H₂ and O₂). The solid phases in the system were pure rutile and a quartz-rutile solid-solution phase. Thermodynamic properties of solid-solution end-members were taken from the Holland and Powell (1998) database (as revised in Thermocalc datafile ds55). The regular binary mixing

model (Anderson and Crerar, 1993) was used, where the end-member activity coefficient is calculated from a single interaction parameter:

$$RT \ln \gamma_{SiO_2} = X_{TiO_2} (1 - X_{SiO_2}) W_{SiO_2-TiO_2} \quad (16)$$

$$RT \ln \gamma_{TiO_2} = X_{SiO_2} (1 - X_{TiO_2}) W_{SiO_2-TiO_2} \quad (17)$$

Here, γ_{SiO_2} and γ_{TiO_2} represent the activity coefficients, and X_{SiO_2} and X_{TiO_2} the mole fractions of SiO_2 and TiO_2 in quartz. $W_{SiO_2-TiO_2}$ is the regular interaction parameter, which is a simple function of temperature and pressure:

$$W_{SiO_2-TiO_2} = a + bT + cP \quad (18)$$

Here, T is the temperature in Kelvin and P is the pressure in bar. The coefficients a , b and c are adjustable parameters in the regression of experimental data (measured Ti concentrations in quartz at given temperature and pressure).

For the first fitting task (A1), a global optimization setup was prepared (Table 8). The chosen global fitting algorithm was *GN_ISRES* (Runarsson and Xin, 2005), and the weighting used was the inverse square of the measured value. Results of fitting tasks are listed in Table 8, and calculations from the model are compared to experimental data in Fig. 8. Optimized values of the a , b and c coefficients obtained from the global optimization run were then used as initial values for the second task (A2) using the BOBYQA (Powell, 1994) local optimization algorithm. Parameter correlation coefficients using the local algorithm are given in Table 9, showing that coefficients a and c are highly correlated (-0.91). Parameter composite scaled sensitivities from the statistical analysis show that the second coefficient (b), which describes the temperature dependence of the interaction parameter, is the least sensitive to the experimental data.

A third fitting task (A3) was produced using the initial values of parameters obtained from the first task (A1), but setting the calculated and measured concentrations of Ti in quartz in the natural logarithm scale. The resulting parameter values are almost identical to the ones obtained from the second fitting task, but their computed errors (2 times the standard deviation of the parameters resulting from 1000 Monte Carlo runs) are half as large as the errors from the second fitting task (Table 8). The change in error values is a consequence of the change in the shape of the minimized function surface due the conversion to logarithmic scale.

Table 8. Results of optimization runs for the Ti-in-quartz solid solution model.

Coefficient	(1)Initial	(1)Final (global)	(2)Final (local)	(2)Error (2 σ)	(2)CSS	(3)Final (local)	(3)Error (2 σ)
<i>a</i>	10000 (1000 – 100000)	60300	60316	2300	35.4	60717	1100
<i>b</i>	-1 (-100 – 100)	-1.168	-1.159	0.568	1.32	-1.577	0.42
<i>c</i>	1 (-100 – 100)	1.791	1.780	0.180	13.5	1.762	0.1

Coefficients *a*, *b* and *c* of the interaction parameter $W = a + bT + cP$. The uncertainty (2 σ) represents the 2 times standard deviation of the parameters from 1000 Monte Carlo runs. Numbers in parentheses represent the parameter bounds during optimization. CSS: composite scaled sensitivities.

(1)Case A1, initial and final values as well as errors for the runs where the global optimization algorithm was applied.

(2)Case A2, final values and errors obtained from the optimization runs using the local optimization algorithm. The final values of the runs with the global optimization algorithm were used as initial guesses for the subsequent runs with the local optimization algorithm.

(3)Case A3, final values and errors obtained from the optimization runs where the measured Ti concentration in quartz was used as $\ln(X_{TiO_2}^{Quartz})$

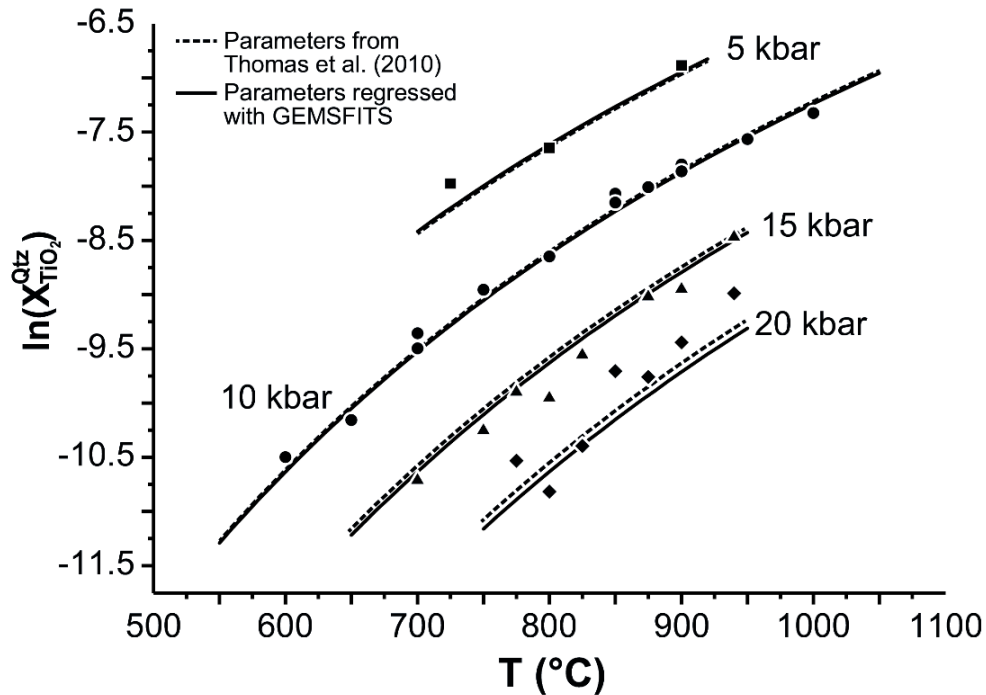


Fig. 8. Plot illustrating the temperature dependence of Ti concentrations in quartz. Filled circles are experimental data points (Wark and Watson, 2008; Thomas et al., 2010) and curves are calculated from thermodynamic solid-solution models.

Table 9. Correlation matrix for the final regressed values of the interaction parameter coefficients ($W = a + bT + cP$) of the Ti in quartz solid solution model.

	<i>a</i>	<i>b</i>	<i>c</i>
<i>a</i>	1.0	0.32	-0.91
<i>b</i>		1.0	-0.02
<i>c</i>			1.0

Thomas et al. (2010) used the following equation to describe their data:

$$RT \ln X_{\text{TiO}_2}^{\text{quartz}} = -a + bT(K) - cP(\text{kbar}) + RT \ln a_{\text{TiO}_2} \quad (19)$$

where R is the universal gas constant, T is temperature in Kelvin, P is pressure in kbar, $\ln X_{\text{TiO}_2}^{\text{quartz}}$ is the mole fraction of TiO_2 in quartz, and a_{TiO_2} is the activity of TiO_2 in the system. Their fitted values for the coefficients (a : -60952 ± 3177 ; b : 1.52 ± 0.39 ; c : 1741 ± 63) agree within uncertainty with those obtained from our regression. This is because eq. (19) can be converted into a form that is equivalent to the regular mixing model (eqs. 16, 17 and 18) using the definition of the activity of end-members:

$$\ln X_{\text{TiO}_2}^{\text{Quartz}} = \ln a_{\text{TiO}_2} - \ln \gamma_{\text{TiO}_2}^{\text{Quartz}} \quad (20)$$

For the binary regular model:

$$RT \ln \gamma_{\text{TiO}_2}^{\text{Quartz}} = (X_{\text{SiO}_2}^{\text{Quartz}})^2 W_{\text{SiO}_2-\text{TiO}_2} \quad (21)$$

When Ti is present at trace element concentrations in quartz, the mole fraction of the SiO_2 end-member can be closely approximated with unity, leading to the simplified expression:

$$\ln \gamma_{\text{TiO}_2}^{\text{Quartz}} = \frac{W_{\text{SiO}_2-\text{TiO}_2}}{RT} \quad (22)$$

This can be substituted into eq. (20) to arrive at a form that is identical to eq. (19).

In an application to field-based geochemical data, the optimized coefficients obtained from the second fitting task (Table 8) were used to determine the temperature of crystallization of quartz (Table 10) in the presence of rutile from 6 natural samples reported by Kidder et al. (2013). For this inverse modeling problem, the Ti concentration data from Kidder et al. (2013) and our own data from quartz veins in low-

grade metamorphic rocks from the Central Swiss Alps (Miron et al., 2013) were added to the TiQ database file. All samples had 300 °C as starting value for the temperature. The parameter marked for fitting in the GEMSFITS task input file was temperature "TK" and in the "DataTarget" section, a nested function was used to compare the measured concentrations of Ti in quartz with the calculated ones. The `gemsfit2` code adjusted the temperatures for each sample independently, until the calculated amount of Ti in solid-solution with quartz was close to the analyzed concentrations. The resulting equilibrium temperatures are in good agreement with the temperatures estimated from other independent geothermometers for the same samples. Other applications similar to the one above could be developed for modeling of trace element partitioning. One can easily use GEMSFITS for adjusting the mixing model parameters for different solid-solutions between major and trace elements in minerals (e.g. Zr in rutile, Zr in titanite, etc.).

Table 10. Application of Ti in quartz geothermometry to two field examples.

Sample	Ti in quartz (ppm)	⁽¹⁾ P (bar)	⁽¹⁾ T (°C)	⁽¹⁾ Error (2σ)	⁽²⁾ T (°C)	⁽²⁾ Error (2σ)	Reference
ms-004	0.19	2600	237	26	234	24	K
ms-34	0.48	3100	276	28	273	22	K
ms-342	0.52	3000	279	28	277	22	K
q-005	0.54	2800	282	28	274	22	K
q-123b	0.30	3100	256	27	257	21	K
q-148j	0.72	3200	295	29	292	23	K
Thusis	1.00	3000*	320*	20*	302	23	M

Temperatures were calculated using analyzed Ti concentrations in quartz, applying the solution model calibration from (Thomas et al., 2010) and this study.

⁽¹⁾Sample median temperatures reported by Kidder et al. (2013) (ref. 'K') and calculated by using TitaniQ thermometer (Thomas et al., 2010) and (*) estimated using mineral geothermometry (Miron et al., 2013). Pressures calculated using a geothermal gradient of 25°Km⁻¹ and the mean temperature values reported by Kidder et al. (2013). Ref. 'M': G.D.Miron, unpublished data.

⁽²⁾Temperatures calculated using inverse modeling with the solution model calibration obtained in this study. The analytical 2σ error has been propagated from the error on the regressed interaction parameters.

6. Discussion

The GEMSFITS code package can adjust separately or simultaneously any GEM input property parameters (standard state Gibbs energies of formation, interaction parameters of thermodynamic mixing models or equations of state, pressure, temperature, input bulk composition), provided that sufficient experimental data are available that can be compared to their computed counterparts. This is a substantial extension compared

to the previous prototype (Hingerl et al., 2014), which had only capabilities for fitting the interaction parameters of aqueous activity models.

The user-defined objective function with one or multiple terms (one term for each type of property) makes it possible to calculate the sum of residuals for any measurable property from experiments performed in different, but related systems. GEMSFITS can therefore fit parameters simultaneously for several chemical systems and many individual experimental data points. The quality of fit can be improved by assigning conditional weights to experimental data points using a range of methods (Table 3). Weights can also be placed on each term of the objective function, thereby giving more weight to selected types of measured data (e.g., placing more weight on solubility data than on volumetric or calorimetric data of aqueous electrolytes). Weighting can also be used to normalize the observed and computed data, because even for a linear regression model it is recommended that all data sets are expressed in the same scale (Motulsky and Christopoulos, 2004). For datasets that cover a large range in parameter space (e.g., experimental solubility data that span several orders of magnitude in concentration), the data should be brought to the same units or at least to similar magnitudes, in order to avoid that the high-magnitude properties would strongly bias the fitting results.

In the `gemsfit2` code, one way of bringing the experimental data to the same magnitude is by using the squared inverse measured value as a weight multiplier. For certain types of data such as dissolved aqueous concentrations, the logarithmic scale is clearly preferable. A potential problem when using the logarithmic scale is that the residuals are asymmetric, i.e., that the absolute value of the positive residual is numerically larger than that of the negative residual (even if the residuals in the non-logarithmic scale are absolutely equal). It is often difficult to decide which experimental points should be treated as outliers. Samples identified as outliers can be skipped from the fitting task, or assigned a low weight (close to zero). A general way of treating outliers is using the GEMSFITS implementation of the Tuckey's Biweight method (Motulsky and Christopoulos, 2004), which gives increasingly less weight to the points further away from the ideal value.

Well constrained parameters for mixing models in non-ideal solution phases are important for accurately calculating activities of solutes or end-members, which is essential for realistic modeling of (geo)chemical systems. These models can have a large number of parameters (e.g., the Pitzer model; Pitzer and Kim, 1974), which imposes high demands on the fitting method. Therefore, it is important to utilize

numerically robust and stable optimization algorithms and to perform thorough statistical analysis of the fitting results. The `gemfit2` code therefore provides a selection of different global and local, gradient based or gradient-free, algorithms that can solve multidimensional non-linear optimization problems. Local optimization algorithms may experience problems when the surface of the function to be minimized is complex and has local minima. In this case, the parameter optimization may converge to incorrect values that represent a local rather than the global minimum. Global algorithms are much slower, but can search for the true global minimum. However, global algorithms are less precise close to the minimum compared to the local algorithms. Thus, a ‘smart’ procedure combines a global optimization run and then uses the resulting parameters as initial guesses with rather close upper and lower bounds in a subsequent local refinement step. Compared to a global algorithm, the parameter search space of a local algorithm is smaller, and this increases the calculation time required to find the optimal solution of the fitting problem.

The standard state molar (molal) Gibbs energy values $\Delta_f G_{298}^0$ of dependent components (e.g., aqueous species, solid-solution end-members, pure condensed substances, gaseous and fluid species) can be adjusted as free (‘F’) parameters, or by application of reaction-type (‘R’) constraints. Complex chemical systems contain many dependent components, which considerably increases the dimensionality of the fitting problem and makes it more difficult for the optimization algorithms to converge. Furthermore, some of the regressed $\Delta_f G_{298}^0$ values may be highly correlated and may not be physically meaningful. For complex systems, the reaction-type constraints can greatly speed up the fitting process for poorly constrained or highly correlated $\Delta_f G_{298}^0$ values, especially if only low-quality or insufficient experimental data are available (see the example on fitting the $\Delta_f G_{298}^0$ values of aqueous aluminum species).

Inverse GEM modeling tasks include thermobarometry calculations and inverse titrations. Thermobarometry finds the temperature and/or pressure based on analyzed compositions of solution phases (e.g., Mg-Fe exchange equilibria between garnet and biotite, Ti concentration in quartz etc.). Mineral solid-solutions record the P - T history of the rocks by continuous adjustment of the element partitioning between coexisting mineral phases, driven by changes in P - T conditions (Spear, 1993; Powell and Holland, 2008). The response and sensitivity of element exchange reactions to changes in temperature and pressure is determined from well-defined laboratory experiments or analytical data for rock samples where independent information on P - T conditions is available (Zhou et al., 1994; Dale et al., 2000; Worley and Powell, 2000).

These data are used to calibrate thermodynamic models that describe the mixing properties of the mineral solid-solution phases involved in the element partition reactions. Thermodynamic models can then be used to estimate the P - T conditions that natural rock samples have experienced.

Inverse titrations involve the iterative adjustment of the bulk composition of the chemical system to match the calculated output properties with their given (experimental) counterparts (e.g., pH, p_e , activities/fugacities of species in gas/fluid phases or in aqueous solution). Commonly, the exact amounts of titrants employed in the experiments to adjust some parameters such as pH are not reported in the publications, but only the measured output parameters (e.g. pH) are provided. For such cases, “nested” objective functions can be defined in GEMSFITS, and parameter optimization can be performed using experimental data such as measured mineral solubility as function of pH (as in the boehmite solubility example), or pH edges for adsorption of aqueous ions on solid surfaces. The `gemsfit2` code then adjusts the amount of titrant through the nested inverse titration functions, until the computed equilibrium pH is in agreement with the measured pH at a prescribed precision. This option has to be treated carefully, because it could result in undesired changes in the ionic strength of aqueous solutions.

In GEMSFITS, the experimental data are currently kept in a NoSQL database as local files within the project folder. The database contents can be exported to JSON text files for backup or further to be uploaded to a MongoDB server, if necessary. Maintenance of the database and editing data records can be straightforwardly performed through the `gfshell12` graphical user interface. Compared to the SQL database that was used in the early prototype (Hingerl et al., 2014), the NoSQL database is much better suited for storing weakly structured data that describe samples with variable experimental conditions. In the NoSQL database, there is no need to know the data structure in advance before creating the database or even inserting new records (documents) into an existing database. The way in which the experiments are stored, handled and selected in GEMSFITS permits to test the fitting of many data combinations from different experimental settings, as well as to remove experimental data sets or single outliers without the need to prepare different experimental data input files for each fitting task. The latter is often required by other fitting tools (Herbelin, 1999; Karamalidis and Dzombak, 2010; Kinniburgh and Cooper, 2011).

A major advantage of GEMSFITS is that both standard statistical and Monte Carlo based methods are available for analysis of the regression results. Monte Carlo methods are essential if the documentation of the analytical errors of the experimental data in the original publications is unavailable, incomplete, or

lacking appropriate consideration of all sources of uncertainty. When using the Monte Carlo method, performing global optimization of large systems, or optimizing a large number of fitting parameters, the `gemsfit2` code can take advantage of the parallelization, which substantially decreases the computing time required to complete a fitting task.

For example, fitting the Ti-in-quartz solid-solution model required to optimize 3 parameters using the global algorithm (ISRES) with 20000 iterations. To complete this task, the program executed 3.2 times faster when parallelized on 4 processor threads compared to one thread (32 compared to 101 seconds). When fitting standard state properties of aqueous Al species (using 4 free parameters and one parameter constrained by the species-forming reaction) with a local algorithm (BOBYQA), the program needed 54 iterations to converge and executed 3.4 times faster when parallelized on 4 processor threads compared to one thread (644 compared to 2220 seconds). The speedup gained from parallelization will become more important in large chemical systems, where the number of fitting parameters and number of experiments will dramatically increase. Producing internally consistent thermodynamic datasets for large chemical systems will involve simultaneous regression of many standard state Gibbs energies of species using thousands of experimental data points. This work will only be feasible when taking advantage of the code parallelization that GEMSFITS offers.

The GEMSFITS codes will be made freely distributable and open-source, as part of the GEM Software collection (<http://gems.web.psi.ch>). This will give other scientists the opportunity to use the codes and the possibility to improve them further. Free (as ‘freedom’) software is of great importance in modern research communities where scientists share their knowledge in a way that others can build upon and use it freely in their research.

7. Outlook

All the new features described above make the GEMSFITS code package a general, flexible, efficient, and user-friendly practical tool for fitting any input parameters of geochemical-thermodynamic models. Future implementations will include extensions to fit the parameters of electrostatic sorption models and mineral dissolution/precipitation kinetic models that are part of the current development version of the GEMS3K codes. When completed, these models will be included into the release version of GEMS and made available to the scientific community. This will increase the range of applications of GEMSFITS to surface adsorption

studies and mineral-aqueous reaction kinetics. Furthermore, GEMSFITS will be extended to be capable of fitting parameters of all equation-of-state models that are implemented into the GEM-Selektor v.3 code. Another long-term goal is to create an experimental database server accessible online via web applications that would be updated and improved using the scientific expertise and resources of different participants. This data can then be easily selected and used with GEMSFITS for optimizing various models for a large number of applications in geochemistry, petrology, chemical engineering, and materials science.

Acknowledgments

This project was supported by ETHIRA grant ETH-19-11-2 from ETH Zürich. Additional funding to DK was provided by Nagra, Wettingen. We thank Pawel Kuczera for helpful discussions concerning parameter fitting approaches and related statistical data analysis.

References

- Anderson, G.M., Crerar, D.A., 1993. *Thermodynamics in Geochemistry: The equilibrium model*. Oxford University Press, 588 p.
- Bénézech, P., Palmer, D.A., Wesolowski, D.J., 2001. Aqueous high-temperature solubility studies. II. The solubility of boehmite at 0.03 m ionic strength as a function of temperature and pH as determined by in situ measurements. *Geochim. Cosmochim. Acta* 65, 2097-2111.
- Brent, R., 1972. *Algorithms for Minimization without Derivatives* Prentice-Hall, Englewood Cliffs, New Jersey, 195 pp.
- Rusk, B.G., Lowers, H.A., Reed, M.H., 2008. Trace elements in hydrothermal quartz: relationships to cathodoluminescence textures and insights into vein formation. *Geology* 36, 547–550.
- Dale, J., Holland, T., and Powell, R., 2000. Hornblende-garnet-plagioclase thermobarometry; a natural assemblage calibration of the thermodynamics of hornblende. *Contrib. Mineral. Petrol.* 140, 353-362.
- DiCiccio, T.J., Efron, B., 1996. Bootstrap confidence intervals (with Discussion). *Stat. Sci.* 11, 189-228.
- Doherty, J. E., Hunt, R.J., 2010, *Approaches to Highly Parameterized Inversion: A Guide to Using PEST for Groundwater-Model Calibration*: US Geological Survey.

- Eldred, M.S., Giunta, A.A., van Bloemen Waanders, B.G., Wojtkiewicz, S.F., Hart, W.E., Alleva, M.P., 2007. DAKOTA, a multilevel parallel object-oriented framework for design optimization, parameter estimation, uncertainty quantification, and sensitivity analysis: Version 4.1. Sandia National Laboratories Albuquerque.
- Gablonsky, J.M., Kelley, C.T., 2001. A locally-biased form of the DIRECT algorithm. *J. Global Optimization* 21, 27-37.
- Helgeson, H.C., Kirkham, D.H., Flowers, G.C., 1981. Theoretical prediction of the thermodynamic behavior of aqueous electrolytes by high pressures and temperatures. IV. Calculation of activity coefficients, osmotic coefficients, and apparent molal and standard and relative partial molal properties to 600 °C and 5kb. *Amer. J. Sci.* 281, 1249-1516.
- Hemingway, B.S., Robie, R.A., Apps, J.A., 1991. Revised values for the thermodynamic properties of boehmite, AlO(OH), and related species and phases in the system Al-H-O. *Amer. Mineral.* 76, 445-457.
- Herbelin, A.L., 1999. FITEQL a computer program for determination of chemical equilibrium constants from experimental data. Oregon State University, Corvallis.
- Hill, M.C., Tiedeman, C.R., 2007. Effective groundwater model calibration: With analysis of data, sensitivities, predictions, and uncertainty.
- Hingerl, F., Kosakowski, G., Wagner, T., Kulik, D., Driesner, T., 2014. GEMSFIT: a generic fitting tool for geochemical activity models: *Computat. Geosci.* 18, 227-242.
- Hoffmann, J., Kräutle, S., Knabner, P., 2012. A general reduction scheme for reactive transport in porous media. *Computat. Geosci.* 16, 1081–1099.
- Holland, T.J.B., Powell, R., 1998. An internally consistent thermodynamic data set for phases of petrological interest. *J. Metam. Geol.* 16, 309-343.
- Huang, R., Audétat, A., 2012. The titanium-in-quartz (TitaniQ) thermobarometer: A critical examination and re-calibration. *Geochim. Cosmochim. Acta* 84, 75-89.
- Johnson, J.W., Oelkers, E.H., Helgeson, H.C., 1992. SUPCRT92: A software package for calculating the standard molal thermodynamic properties of minerals, gases, aqueous species, and reactions from 1 to 5000 bars and 0°C to 1000°C. *Comput. Geosci.* 18, 890–947.

- Jones, D.R., Perttunen, C.D., Stuckmann, B.E., 1993. Lipschitzian optimization without the lipschitz constant. *J. Optim. Theory Appl.* 79, 157.
- Kaelo, P., Ali, M.M., 2006. Some variants of the controlled random search algorithm for global optimization. *J. Optim. Theory Appl.* 130, 253-264.
- Karamalidis, A.K., Dzombak, D.A., 2010. Data compilation and treatment methods. *Surface Complexation Modeling*, John Wiley & Sons, 45-57.
- Karpov, I.K., Chudnenko, K.V., Kulik, D.A., 1997. Modeling chemical mass transfer in geochemical processes; thermodynamic relations, conditions of equilibria and numerical algorithms. *Amer. J. Sci.* 297, 767-806.
- Kestin, J., Sengers, J.V., Kamgar-Parsi, B., Levelt-Sengers, J.M., 1984. Thermophysical properties of fluid H₂O. *J. Phys. Chem. Ref. Data* 13, 175-183.
- Kidder, S., Avouac, J.P., Chan, Y.C., 2013. Application of titanium-in-quartz thermobarometry to greenschist facies veins and recrystallized quartzites in the Hsüehshan range, Taiwan. *Solid Earth* 4, 1-21.
- Kinniburgh, D., Cooper, D., 2011. PhreePlot: Creating graphical output with PHREEQC.
- Kraft, D., 1994. Algorithm 733: TOMP–Fortran modules for optimal control calculations. *ACM Transact. Mathem. Software* 20, 262-281.
- Kulik, D.A., Wagner, T., Dmytrieva, S.V., Kosakowski, G., Hingerl, F., Chudnenko, K.V., Berner, U., 2013. GEM-Selektor geochemical modeling package: revised algorithm and GEMS3K numerical kernel for coupled simulation codes. *Computat. Geosci.* 17, 1-24.
- Leal A.M.M., Blunt M, LaForce T., 2014. Efficient chemical equilibrium calculations for geochemical speciation and reactive transport modelling, *Geochim. Cosmochim. Acta* 131, 301-322.
- Manning, C.E., 2006. Mobilizing aluminum in crustal and mantle fluids. *J. Geochem. Explor.* 89, 251-253.
- MathWorks, 2012. MATLAB and Statistics Toolbox Release 2012b, The MathWorks Inc., Natick, Massachusetts, United States.
- Miron, G., Wagner, T., Wälle, M., Heinrich, C.A., 2013. Major and trace-element composition and pressure–temperature evolution of rock-buffered fluids in low-grade accretionary-wedge metasediments, Central Alps. *Contrib. Mineral. Petrol.* 165, 981-1008.

- Mookherjee, M., Keppler, H., Manning, C.E., 2014. Aluminum speciation in aqueous fluids at deep crustal pressure and temperature. *Geochim. Cosmochim. Acta* 133, 128-141.
- Motulsky, H., Christopoulos, A., 2004. Fitting models to biological data using linear and non-linear regression: a practical guide to curve fitting. Oxford University Press, Oxford, 352 p.
- Oelkers, E.H., Helgeson, H.C., 1990. Triple-ion anions and polynuclear complexing in supercritical electrolyte solutions. *Geochim. Cosmochim. Acta* 54, 727-738.
- Ostapenko, G.T., Gamarnik, M.Y., Gorogotskaya, L.I., Kuznetsov, G.V., Tarashchan, A.N., Timoshkova, L.P., 1987. Isomorphism of titanium substitution for silicon in quartz: experimental data. *Mineral Zh.* 9, 30-40.
- Palmer, D.A., Benezeth, P., Wesolowski, D.J., 2001. Aqueous high-temperature solubility studies. I. The solubility of boehmite as functions of ionic strength (to 5 molal, NaCl), temperature (100-290°C), and pH as determined by in situ measurements. *Geochim. Cosmochim. Acta*, 65, 2081-2095.
- Palmer, D.A., Wesolowski, D.J., 1993. Aluminum speciation and equilibria in aqueous solution: III. Potentiometric determination of the first hydrolysis constant of aluminum(III) in sodium chloride solutions to 125°C. *Geochim. Cosmochim. Acta* 57, 2929-2938.
- Parkhurst, D.L., Appelo, C.A.J., 2013. Description of input and examples for PHREEQC version 3: a computer program for speciation, batch-reaction, one-dimensional transport, and inverse geochemical calculations. U.S. Geological Survey Techniques and Methods 6, 497.
- Pitzer, K.S., Kim, J.J., 1974. Thermodynamics of electrolytes. IV. Activity and osmotic coefficients for mixed electrolytes. *J. Amer. Chem. Soc.* 96, 5701-5707.
- Plantenga, T.D., 2009. HOPSPACK2.0 user manual: Version 2.0.2. Sandia Technical Report 2009-6265.
- Plugge, E., Hawkins, T., Membrey, P., 2010. The definitive guide to MongoDB: The NoSQL database for cloud and desktop computing, APress, 328 p.
- Poeter, E.P., Hill, M.C., 1998. Documentation 722 of UCODE: A computer code for universal inverse modeling. U.S. Geological Survey, Denver.
- Pokrovski, G.S., Schott, J., Salvi, S., Gout, R., Kubicki, J.D., 1998. Structure and stability of aluminum-silica complexes in neutral to basic solutions. Experimental study and molecular orbital calculations. *Min. Mag.*, 62A, 1194-1195.

- Powell, M.J.D., 1994. A direct search optimization method that models the objective and constraint functions by linear interpolation. Proceedings of the Sixth Workshop on Optimization and Numerical Analysis, Oaxaca, Mexico, 51.
- Powell, M.J.D., 2004. The NEWUOA software for unconstrained optimization without derivatives. Proceedings of the 40th Workshop on Large Scale Nonlinear Optimization (Erice, Italy, 2004).
- Powell, R., Holland, T.J.B., 2008. On thermobarometry. *J. Metam. Geol.* 26, 155-179.
- Reed, M.H., 1982. Calculation of multicomponent chemical equilibria and reaction processes in systems involving minerals, gases and an aqueous phase. *Geochim. Cosmochim. Acta* 46, 513-528.
- Rinnooy, A.H.G.K, Timmer, G.T., 1987. Stochastic global optimization methods. *Mathem. Progr.* 39, 27-78.
- Rowan, T., 1990. Functional Stability Analysis of Numerical Algorithms. Ph.D. thesis, Department of Computer Sciences, University of Texas at Austin.
- Runarsson, T.P., Xin, Y., 2005. Search biases in constrained evolutionary optimization. *IEEE Transact.* 35, 233-243.
- Silva Santos, C.H., Goncalves, M.S., Hernandez-Figueroa, H.E., 2010. Designing novel photonic devices by bio-inspired computing. *IEEE Photonics Tech. Lett.* 22, 1177-1179.
- Smith, V.C., Shane, P., Nairn, I., 2010. Insights into silicic melt generation using plagioclase, quartz and melt inclusions from the caldera-forming Rotoiti eruption, Taupo volcanic zone, New Zealand. *Contrib. Mineral. Petrol.* 160, 951-971.
- Spear, F.S., 1993. Metamorphic phase equilibria and Pressure-Temperature-Time paths. Mineralogical Society of America, 799 pp.
- Steeffel, C.I., Lasaga, A.C., 1994. A coupled model for transport of multiple chemical species and kinetic precipitation/dissolution reactions with application to reactive flow in single phase hydrothermal systems. *Amer. J. Sci.* 294, 529-592.
- Steeffel, C.I., DePaolo, D.J., Lichtner, P.C., 2005. Reactive transport modeling: An essential tool and a new research approach for the Earth sciences. *Earth Planet. Sci. Lett.* 240, 539-558.
- Shock, E.L., Helgeson, H.C., 1988. Calculation of the thermodynamic and transport properties of aqueous species at high pressures and temperatures: Correlation algorithms for ionic species and equation of state predictions to 5 kb and 1000°C. *Geochim. Cosmochim. Acta* 52, 2009-2036.

- Shock, E.L., Oelkers, E.H., Johnson, J.W., Sverjensky, D.A., Helgeson, H.C., 1992. Calculation of the thermodynamic properties of aqueous species at high pressures and temperatures. Effective electrostatic radii, dissociation constants and standard partial molal properties to 1000 °C and 5 kbar. *J. Chem. Soc. Farad. Trans.* 88, 803-826.
- Shock, E. L., Sassani, D. C., Willis, M., and Sverjensky, D. A., 1997. Inorganic species in geologic fluids: Correlations among standard molal thermodynamic properties of aqueous ions and hydroxide complexes. *Geochim. Cosmochim. Acta* 61, 907-950.
- Svanberg, K., 2002. A class of globally convergent optimization methods based on conservative convex separable approximations. *SIAM J. Optim.* 12, 555-573.
- Sverjensky, D., Shock, E., Helgeson, H., 1997. Prediction of the thermodynamic properties of aqueous metal complexes to 1000 C and 5 kb. *Geochim. Cosmochim. Acta* 61, 1359-1412.
- Tagirov, B.R., Zotov, A.V., Akinfiev, N.N., 1997. Experimental study of dissociation of HCl from 350 to 500°C and from 500 to 2500 bars: Thermodynamic properties of HCl⁰(aq). *Geochim. Cosmochim. Acta* 61, 4267-4280.
- Tagirov, B., Schott, J., 2001. Aluminum speciation in crustal fluids revisited. *Geochim. Cosmochim. Acta* 65, 3965-3992.
- Tanger, J.C., Helgeson, H.C., 1988. Calculation of the thermodynamic and transport properties of aqueous species at high pressures and temperatures: Revised equations of state for the standard partial molal properties of ions and electrolytes. *Amer. J. Sci.* 288, 19-98.
- Thomas, J.B., Watson, E.B., Spear, F.S., Shemella, P.T., Nayak, S.K., Lanzirrotti, A., 2010. TitaniQ under pressure: the effect of pressure and temperature on the solubility of Ti in quartz. *Contrib. Mineral. Petrol.* 160, 743-759.
- Verdes, G., Gout, R., Castet, S., 1992. Thermodynamic properties of the aluminate ion and of bayerite, boehmite, diaspore and gibbsite. *Eur. J. Mineral.* 4, 767-792.
- Vlcek, J., Luksan, L., 2006. Shifted limited-memory variable metric methods for large-scale unconstrained minimization. *J. Computat. Appl. Math.* 186, 365-390.
- Wagner, T., Kulik, D.A., Hingerl, F.F., Dmytrieva, S.V., 2012. GEM-Selektor geochemical modeling package: TSolMod library and data interface for multicomponent phase models. *Can. Mineral.* 50, 1173-1195.

- Westall, J.C., Zachary, J.L., Morel, F.M.M., 1976. MINEQL: A computer program for the calculation of chemical equilibrium composition of aqueous systems. Inst. of Technol., Cambridge.
- Wilson, C., Seward, T., Allan, A., Charlier, B., Bello, L., 2012. A comment on: 'TitaniQ under pressure: the effect of pressure and temperature on the solubility of Ti in quartz', by Jay B. Thomas, E. Bruce Watson, Frank S. Spear, Philip T. Shemella, Saroj K. Nayak and Antonio Lanzirrotti. *Contrib. Mineral. Petrol.* 164, 359.
- Worley, B., Powell, R., 2000. High-precision relative thermobarometry; theory and a worked example. *J. Metam. Geol.* 18, 91-101.
- Xu, T., Spycher, N., Sonnenthal, E., Zhang, G., Zheng, L., Pruess, K., 2011. TOUGHREACT Version 2.0: A simulator for subsurface reactive transport under non-isothermal multiphase flow conditions. *Comp. Geosci.* 37, 763-774.
- Zhang, F., Yeh, G.T., Parker, J.C. (Eds.), 2012. *Groundwater Reactive Transport Models*. Behtham Publishers, 254 p.
- Zhou, T., Dong, G., Phillips, G. N., 1994. Chemographic analysis of assemblages involving pyrophyllite, chloritoid, chlorite, kaolinite, kyanite, quartz; application to metapelites in the Witwatersrand goldfields, South Africa. *J. Metam. Geol.* 12, 655-666.