



HAL
open science

Comparison of linear solvers for equilibrium geochemistry computations

Hela Machat, Jérôme Carrayrou

► **To cite this version:**

Hela Machat, Jérôme Carrayrou. Comparison of linear solvers for equilibrium geochemistry computations. *Computational Geosciences*, 2017, 21 (1), pp.131-150. 10.1007/s10596-016-9600-5. hal-03321650

HAL Id: hal-03321650

<https://hal.science/hal-03321650>

Submitted on 17 Aug 2021

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

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

2 **Comparison of linear solvers for equilibrium geochemistry**
3 **computations**4 **Hela Machat^{1,2} · Jérôme Carrayrou¹**

5 Received: 11 February 2015 / Accepted: 24 October 2016

6 © Springer International Publishing Switzerland 2016

7 **Abstract** Equilibrium chemistry computations and reactive
8 transport modelling require the intensive use of a linear
9 solver under very specific conditions. The systems to be
10 solved are small or very small (4×4 to 20×20 , occasionally
11 larger) and are very ill-conditioned (condition number up to
12 10^{100}). These specific conditions have never been investi-
13 gated in terms of the robustness, accuracy, and efficiency of
14 the linear solver. In this work, we present the specificity of
15 the linear system to be solved. Several direct and iterative
16 solvers are compared using a panel of chemical systems,
17 including or excluding the formation of mineral species. We
18 show that direct and iterative solvers can be used for these
19 problems and propose computational keys to improve the
20 chemical solvers.

21 **Keywords** Geochemical modelling · Instantaneous
22 equilibrium chemistry · Linear system inversion · Linear
23 solver · Small matrix · Ill-conditioned matrix ·
24 Newton-Raphson algorithm

25 **1 Introduction**

26 The problem of groundwater management is receiving
27 increasing attention, and many tools have been developed

to address this issue. One of these tools, reactive trans- 28
port models, was first limited to laboratory experiments 29
and was then extended to field problem comprehension. In 30
recent decades, reactive transport models have increased in 31
complexity and efficiency, and they are now used in many 32
fields. Reactive transport models have been used to study 33
the transport of contaminants, such as heavy metals [1, 2] 34
and radioelements [3–5]. Because of the increasing inter- 35
est in questions related to climate change, many studies 36
on reactive transport have been conducted to examine the 37
possibility of geologic CO_2 sequestration [6–10]. 38

Under the wide variety of models and cases lies a com- 39
mon mathematical description [11–13]. Transport is usually 40
described by an advection-dispersion equation, and the 41
chemistry is formulated under thermodynamic equilibrium. 42
A widely used approach to solve these reactive transport 43
problems is the operator splitting approach [14]. Using this 44
approach, the transport and chemical operators are solved 45
separately at each time step and iteratively for some for- 46
mulations. As a consequence, the chemistry operator has to 47
be solved at least once per mesh cell per time step. This 48
is one reason for the high computational cost of reactive 49
transport modelling. Some authors have reported that 80 50
to 90 % of the computation time is dedicated to chem- 51
ical computation. Many studies have been conducted to 52
reduce the computation time required by reactive transport 53
modelling [15]. Some works have explored paralleliza- 54
tion [16], while others have focused on the methods used 55
to solve the transport operator. Nevertheless, improving 56
the resolution of the chemistry operator has been iden- 57
tified as a key point. Some authors have attempted to 58
improve the classic Newton-Raphson method [17], while 59
others have tested other methods, such as Newton-Krylov 60
[16, 18]. 61

✉ Jérôme Carrayrou
jerome.carrayrou@unistra.fr

¹ CNRS, ENGEES, LHyGeS UMR 7517, Université de
Strasbourg, 67000 Strasbourg, France

² Ecole Supérieure des Ingénieurs de l'Équipement Rural de
Medjez el Bab, University of Jendouba, Jendouba, Tunisia

In this work, we focus on a specific element of the problem, improving the resolution of the linearized system provided by the Newton-Raphson method. Looking to numerical methods to solve linear systems is not currently a common practice. Indeed, these methods are actually well known [19–23], and all mathematical packages for scientific computation propose several routines for this task. The motivation of this work comes from the specificity of linear systems that have to be solved for equilibrium chemistry computations. Classic tests for the resolution of linear systems [24–30] are performed using systems provided by finite element or finite volume discretization, leading to matrices that are large (at least 10,000 unknowns) and sparse. Moreover, even when ill-conditioned systems have been studied [25, 30, 31], the conditioning of the matrix coming from the chemical system is specific, as underlined by Hoffmann et al. [32]. For example, Soleymani [33] worked with an ill-conditioned system constructed from 10×10 to 20×20 Hilbert matrices. The condition numbers then range from 3.5×10^{13} to 6.2×10^{28} . In this work, we present

chemical tests leading to a 7×7 matrix with a condition number of approximately 10^{180} .

We expect to find a method to increase the efficiency of a speciation or reactive transport code. Several properties are required for such a method:

- (i) This method should be fast, as the linear system will be solved very often. In the case of reactive transport modelling, the system will be solved at least once per mesh cell per time step.
- (ii) The method should be very robust. It should be able to solve the linear system even if it is very poorly conditioned. Because the resolution of the linear system is only part of an iterative Newton step, an accurate solution is not absolutely needed. Thus, some advanced codes (e.g. Linear Algebra Package (LAPACK) routine) that check the accuracy of the solution and return an error flag instead of an inaccurate solution are, in this work, less robust than the more rustic routines.
- (iii) The method should be able to detect failure and return an error flag to the main program so that a recovery procedure can be initiated. In the case of reactive transport modelling, this procedure could involve rejecting the current time step and recomputing with a smaller one.
- (iv) In the initial analysis, the precision of the method is not the key point. Because the linear system resolution is only a part of the Newton-Raphson iterative procedure, *reasonable* error is acceptable for the linear system inversion. If this error is too large, it will slow the convergence speed for the Newton-Raphson method and decrease the efficiency of the reactive transport code. In this work, errors are estimated

by comparing the calculated solution to a reference solution. 114
115

Because we utilize a markedly small matrix, we did not test parallelization. All the computations were performed on a PC running Windows with 64-bit Fortran 95. Real variables are defined as double-precision real. We prefer double-precision computations because all the chemical codes are, to the best of our knowledge, written as double-precision real and because quadruple-precision computation is much more time consuming. Nevertheless, we have tested one method using quadruple-precision real to determine whether this development could be useful. Reference solutions are also computed using quadruple precision. 116
117
118
119
120
121
122
123
124
125
126

We first present the formulation of the equations describing equilibrium reactions and how they are solved using the Newton-Raphson method. This point defines the Jacobian linear system, which is the object of this work. A second part is devoted to the presentation of the chemical tests and the numeric procedures used to perform them. Next, we propose a detailed analysis of the structure and properties

of the Jacobian matrix. The selected linear solvers are then presented and tested, and the results are compared and

discussed. Based on this analysis, we propose an algorithm to optimize the chemical computation in terms of robustness, accuracy, and efficiency. This algorithm is evaluated on the most selective test. By expanding the limits of the currently used methods, we believe that our new algorithm will contribute to enlarging the field of application of reactive transport modelling. As a conclusion, we underline the main advances of this work, the new perspectives and the remaining obstacles. 133
134
135
136
137
138
139
140
141
142
143
144

2 Material and methods 145

2.1 Geochemical modelling 146

One efficient formulation for the computation of thermodynamic equilibrium is based on the tableau concept, referred to as Morel's table [34, 35]. N_X components (X_j) are chosen from the N_C species (C_i) and are used to write the formation of each species as a combination of the components. The mass action law for the formation of the C_i species is written with the equilibrium constant (K_i) and the stoichiometric coefficients ($a_{i,k}$) for each component (X_k) 147
148
149
150
151
152
153
154

$$\{C_i\} = K_i \prod_{k=1}^{N_X} \{X_k\}^{a_{i,k}} \quad (1)$$

where $\{C_i\}$ and $\{X_k\}$ are the activities of species C_i and component X_k , respectively. In this work, we define X_j as a subset of C_i ; then, N_X is N_C minus the number of reactions. 155
156
157

158 If N_{CP} -precipitated species (Cp_i) are taken into account,
 159 the mass action law for the precipitation of Cp_i is written
 160 with the precipitation constant (Kp_i) and the stoichiomet-
 161 ric coefficients ($ap_{i,k}$). The saturation index (SI) of Cp_i is
 162 equal to its activity, which is unity for a pure solid phase

$$SI_i = Kp_i \prod_{k=1}^{N_x} \{X_k\}^{ap_{i,k}} = 1 \quad (2)$$

163 The conservation of the total concentration [T_j] of the j th
 164 component in the system is then written as

$$T_j = \sum_{i=1}^{N_C} a_{i,j} \cdot [C_i] + \sum_{i=1}^{N_{CP}} ap_{i,j} \cdot [Cp_i] \quad (3)$$

165 where $[C_i]$ is the concentration of species (C_i) and $[Cp_i]$ is
 166 the amount of precipitated species (Cp_i) per liquid volume
 167 unit.

168 A classic algorithm [17, 36–41] to describe mineral pre-
 169 cipitation or dissolution makes an a priori hypothesis about
 170 the existence or non-existence of minerals. In this work, we
 171 assume that this hypothesis is proposed. The relationships

$$Y_j = -T_j + \sum_{i=1}^{N_C} a_{i,j} \cdot \prod_{k=1}^{N_x} \left(\frac{K_j}{N_x} \cdot X_k \right)^{a_{i,k}} + \sum_{i=1}^{N_{CP}} ap_{i,j} \cdot [Cp_i] \quad \text{for } j = 1 \text{ to } N_x \quad (6)$$

$$Y_{j=N_x+i} = -1 + Kp_i \cdot \prod_{k=1}^{N_x} (\gamma_k [X_k])^{ap_{i,k}} \quad \text{for } i = 1 \text{ to } N_{CP}$$

185 Using this method, it is possible to include many chemi-
 186 cal phenomena, including activity corrections, sorption on a
 187 surface using different means (such as ion exchange or sur-
 188 face complexation), and dissolution of gaseous compounds.

172 between the activity and concentration are given by activity
 173 coefficients (γ_i) calculated using specific models (Davies,
 174 Debye-Hückel, etc.)

$$\{C_i\} = \gamma_i [C_i] \quad \text{and} \quad \{X_j\} = \gamma_j X_j \quad (4)$$

175 By substituting the mass action law (1) into the mass conser-
 176 vation equation (3), the following relationship, which only
 177 depends on the components and the precipitated species
 178 concentrations, is obtained:

$$T_j = \sum_{i=1}^{N_C} a_{i,j} \cdot \prod_{k=1}^{N_x} (\gamma_k [X_k])^{a_{i,k}} + \sum_{i=1}^{N_{CP}} ap_{i,j} \cdot [Cp_i] \quad (5)$$

179 Combining Eqs. 2 and 5 leads to a set of ($N_x N_{CP}$) non-
 180 linear algebraic equations, which can be numerically solved
 181 through iterative methods. The concentrations of component
 182 $[X_k]$ and precipitated species $[Cp_i]$ at equilibrium are then
 183 determined when the ($N_x N_{CP}$) objective functions (Y_j)
 184 are zero.

185 According to the criteria typically used for this method
 186 [17, 34, 40, 42], the convergence of the Newton-Raphson
 187 method is not checked with respect to the norm of the
 188 objective function $\|Y\|$, but the relative error defined as

$$NR_{\text{relative error}} = \max_j \left| \frac{Y_j}{T_j + \sum_{i=1}^{N_{CP}} ap_{i,j} [Cp_i]} \right|, \quad Y_j: j = N_x + 1, N_x + N_{CP} \leq \epsilon_{N-R} \quad \text{with } \epsilon_{N-R} = 10^{-12} \quad (7)$$

194 The value of the convergence criterion ($\epsilon_{N-R} = 10^{-12}$) is
 195 198
 199

200 formulation has some weaknesses that are explained later
 201 Q4
 set according to usual practice.

2.2 The Newton-Raphson method

The historical approach [12, 34, 37, 40, 42–47] involves the resolution

of the system (6) with the Newton-Raphson method using $[X_k]$ and $[Cp_i]$ as primary unknowns. This

Comput Geosci (see 3.1). Section

202 However, many authors [18, 32, 38, 39, 48] have proposed an alternative approach. Instead of using the component concentrations X_j as the primary variables, they use the logarithm of the component activities ($\ln X_j$). According to this convention, the objective functions defined by Eq. 8 become conservation equations

203 204 205 206 } $\xi = 206$

207 208 -

Q3

$$Y_j = -T_j + \sum_{i=1}^{N_c} a_{i,j} \cdot \frac{K_i}{\gamma_i} \cdot \exp \left(\sum_{k=1}^{N_x} a_{i,k} \cdot \xi_k \right) + \sum_{i=1}^{N_{CP}} a_{p_{i,j}} \cdot [Cp_i] \text{ for } j = 1 \text{ to } N_x \quad (8)$$

209 In the case of the objective function describing precipitation,
 210 it is more interesting to rewrite the mass action law (2) in
 211 log form and then define the objective function

$$Y_{N_x+i} = \ln(SI_i) = \ln(Kp_i) + \sum_{k=1}^{N_x} a_{p_i,k} \cdot \xi_k \text{ for } i = 1 \text{ to } N_{CP} \quad (9)$$

212 Equations 8 and 9 are solved at the n th iteration with the
 213 Jacobian matrix (Z^n) of the objective functions

$$Z^n_{j,k} = \frac{\partial Y^n_j}{\partial \xi_k} \quad j = 1, N_x + N_{CP} \quad k = 1, N_x \quad (10)$$

$$Z^n_{j,k} = \frac{\partial Y^n_j}{\partial C_{p_{k-N_x}}} \quad j = 1, N_x + N_{CP} \quad k = N_x + 1, N_x + N_{CP}$$

214 Z^n can be calculated in two ways.

215 (i) Using an analytical computation, we obtain the $(N_x +$
 216 $N_{CP}) \times (N_x + N_{CP})$ values of Z^n by

$$Z^n_{j,k} = \sum_{i=1}^{N_x} a_{i,j} \cdot a_{i,k} \cdot [C_i]^{n-1} \quad j = 1, N_x \quad k = 1, N_x \quad (11)$$

$$Z^n_{j,k} = a_{p_{k-N_x},j} \quad j = N_x + 1, N_x + N_{CP} \quad k = N_x + 1, N_x + N_{CP}$$

$$Z^n_{j,k} = \sum_{x=1}^{N_x+1} \dots \quad j = 1, N_x + N_{CP} \quad k = 1, N_x + N_{CP}$$

217 229

218 230

219

220

221

222

223

224

225

226

227

228

2.3 Chemical test cases

231

We choose chemical test cases with various numbers of components. Some of these chemical systems allow the formation of mineral species. Although it is not realistic from a chemical point of view, we test them without minerals and with the maximal possible number of minerals to obtain the largest matrix size. Appendix 1 presents the stoichiometric coefficients, equilibrium constants, and concentrations for these tests.

(i) The *gallic acid* test case was presented by Brasard and Bodurtha [49]. It has been recognized as a challenging test for speciation computation [17] (see Appendix 1 (1)).

(ii) The *Valocchi* test is from Valocchi et al. [11]. It involves calcium and magnesium ion exchange (see Appendix 1 (2)).

(iii) The *pyrite* test case describes the dissolution of a pyrite rock in pure water. It has been used to test speciation algorithms [17]. Because it involves redox reactions, the stoichiometric coefficients cover a wide range, and the equilibrium constants vary over several orders of magnitude. This test is used under Appendix 1 (3).

(iv) The *MoMaS easy* test is the chemical system used for the reactive transport benchmark of MoMaS at the easy level [50]. It has been specifically developed to magnify numerical difficulties in a small system (see Appendix 1 (4)).

(v) The *Morel-Morgan* test is the first large chemical system

$$k = N_x + 1, N_x + N_{CP}$$

Even if the activity coefficients depend on the component concentrations, they are assumed to be constant during the Newton-Raphson procedure. These activity coefficients are usually actualized by a fixed-point algorithm at each Newton-Raphson loop.

The progress step of the method ($\Delta \xi^n, \Delta C_p^n$) is achieved by assuming that the objective function Y^{n+1} in Eq. 12 is equal to zero at the $(n + 1)$ th iteration. This produces the key equation of this article, the linear system (12), which must be solved to obtain the progress step

$$Z^n \cdot \Delta \xi^n, \Delta C_p^n = Y^{n+1} - Y^n = -Y^n \quad (12)$$

This system yields the values of the component activities and precipitate concentrations at the $(n + 1)$ th iteration

$$\xi^{n+1} = \xi^n + \Delta \xi^n$$

$$[C_p]^{n+1} = [C_p]^n + \Delta C_p^n \quad (13)$$

To simplify the notations, ξ is used to denote the full vector of unknowns,

including mineral Cp if present.

Comput Geosci reported in the computational literature.

Comput Geosci

was used by F. Morel and M. Morgan in 1972 to present the capacities of the computational method they had just developed (and which we still use today). This test includes 52 components (H^+ , 20 metals, and 31 ligands), leading to 781 aqueous species (see Appendix 1 (5)).

(vi) The *MoMaS medium* test is the chemical system for the medium level of the MoMaS reactive transport benchmark [50] (see Appendix 1 (6)).

(vii) The *Fe-Cr* test is an additional redox test that describes the redox reactions between iron and chromium. These types of reactions occur when iron reactive barriers are used to treat chromium-contaminated sites [51, 52]. In this case, we consider only the aqueous phase without minerals (see Appendix 1 (7)).

(viii) The *pyrite mineral* test describes the dissolution of a pyrite rock in pure water. We assume that three possible mineral phases are present (see Appendix 1 (8)).

282 (ix) The *MoMaS hard* test is the equilibrium part of the
 283 chemical system described in the hard level of the
 284 MoMaS reactive transport benchmark. It allows for
 285 the formation of two mineral species (see Appendix
 286 1 (9)).

287 (x) The *Fe-Cr mineral* test describes the redox reaction
 288 between iron and chromium. We assume the formation
 289 of three different mineral phases (see Appendix
 290 1 (10)).

291 **2.4 Test procedure**

292 Equation 11 shows that we can obtain multiple linear sys-
 293 tems from one chemical problem by changing the activity
 294 values of the components. For each chemical system, we
 295 select three components and vary their values over a wide
 296 range. The concentrations of all minerals are arbitrarily set
 297 to 10^{-3} mol L⁻¹. The activity of component H⁺ is varied
 298 from 10^{-12} to 10^{-2} mol L⁻¹ (pH = 12 to pH = 2), while
 299 that of component e⁻ is varied from 10^{-19} to 10^{12} , corre-
 300 sponding to Eh = 0.7 to 1.1 V computed using Eq. 14 at 25
 301 °C

$$Eh = \ln e^{-\frac{RT}{F}} \quad (14)$$

302 where T is the temperature (Kelvin), R is the gas constant
 303 (8.314 J K mol⁻¹), and F is the Faraday constant (96,487
 304 C mol⁻¹). This range of electrical potential corresponds to
 305 the stability of water at pH values between 2 and 12. For the
 306 O₂ component, it is not possible to cover the same poten-
 307 tial range as e⁻ because of the computation of the reference
 308 solution. The activity is varied from 10^{-70} to 10^4 , as com-

309 puted using Eq. 15 at 25 °C with $E^0 = 1.23$ V and pH
 310 varying from 2 to 12. The potential is then varied from -0.5

311 to 1.1 V

$$Eh = E^0 + \frac{1}{4} \frac{RT}{F} \times \ln \frac{\{O_2\} \{H^+\}^4}{\{H_2O\}} \quad (15)$$

312 The activities of the other components vary from 10^{-12} to
 313 10^{-1} mol L⁻¹. For each of the three selected components,
 314 we compute 30 values equally distributed on a log scale over

315 319
 316 320
 317 321
 318 322

The condition number of Z is defined [23] as the product of 323
 the norm of the matrix per the norm of the inverse matrix 324
 (17) 325

$$\text{cond}(Z) = \|Z\|_1 \times \|Z^{-1}\|_1 \quad (17)$$

To test the numerical methods, we first evaluate the compu- 326
 tation time (CPU time) required to solve the linear system. 327
 Because we work with a very small matrix, the computa- 328
 tions are very fast and we run the same calculation several 329
 times to obtain a total computing time of approximately 1 330

s. The *CPU time* is given in this work in units of seconds 331
 per computation (by dividing the total computing time by 332
 the number of runs). According to this method, the global 333
 computing time for one test case is approximately 6 days. 334

Many numerical methods, including a *failure indicator*, 335
 which indicates the success or failure of the resolution, have 336
 been developed. If needed, we include a failure indicator. 337
 As *failure*, we include the *crash* of the method, underflow 338
 or overflow, non-convergence within the maximum number 339
 of iterations (for iterative methods), or excessive inaccu- 340
 racy for some advanced methods (LAPACK routines) that 341
 estimate the accuracy of the proposed solution. 342

Solving a linear system (13) using a numerical method 343
 produces an approximate solution ($d\xi_{\text{method}}$), and the ref- 344

erence method gives ($d\xi_{\text{ref}}$) with accuracy on the same 345
 order as the roundoff error. To evaluate the accuracy of the 346
 approximate solution, two quantities can be calculated: 347

1. The *relative error on the norm*, Err_{Norm}, is obtained by 348
 computing the norm of the approximate and reference 349
 solution (18) 350

$$\text{Err}_{\text{Norm}} = \frac{\|d\xi_{\text{method}}\| - \|d\xi_{\text{ref}}\|}{\|d\xi_{\text{ref}}\|} \quad (18)$$

1. The error on the direction is given by *angle_{method}*, the 351
 angle (degrees) between the reference and the approx- 352
 imate solution calculated using the scalar product of 353

these two vectors 354

$$\text{angle}_{\text{method}} = \frac{360}{2\pi} \text{Arc cos} \frac{d\xi_{\text{method}} \cdot d\xi_{\text{ref}}}{\|d\xi_{\text{method}}\| \cdot \|d\xi_{\text{ref}}\|} \quad (19)$$

the chosen range, leading to 29,791 different linear systems for each
 chemical test case. For each of these 29,791 tests, we make only one linear
 solver (or one Newton step) (exceptin the last section, Section 4, where the
 iterative Newton method is performed to solve the non-linear system given

by ^{Comput Geosci} All of these quantities, namely the failure indicator, relative error on the norm, angle_{method}, and CPU time, are calculated for the 29,791 linear systems built from each chemical test case for all the tested methods. This enormous amount of data is aggregated in two ways:

norm used in this work is the norm, defined as [23]

(i) For each chemical system and each method, we compute the mean of each quantity.

(ii) For each chemical system and each method, the interval of the condition number is discretized into 100

$$\|Z\|_1 = \max_{i,j} Z_{i,j} \quad (16)$$

regular subintervals. For each subinterval, we compute the mean of each quantity.

$i=1$

366 **2.5 Reference solution**

367 Because of the very high condition numbers, it is not possible to directly obtain an exact solution. We equilibrate the rows and columns of the Jacobian matrices to reduce their condition number using the iterative algorithm proposed by Knight et al. [53] because it preserves the symmetry of the Jacobian matrix.

373 $\tilde{\mathbf{Z}}^k$ be the equilibrated Jacobian matrix at iteration k , with $\tilde{\mathbf{Z}} = \mathbf{Z}$.
 374 These authors defined r_i^k as the vector formed by the i th
 375 k

376 row of $\tilde{\mathbf{Z}}$ and c_i^k as the vector formed by the i th column. The preconditioning matrices \mathbf{R}^k and \mathbf{C}^k are then defined by

$$377 \mathbf{R}^k = \text{diag} \left(\begin{matrix} 1 \\ \vdots \\ r_i^k \\ \vdots \\ \infty \end{matrix} \right)_{i=1, Nx+NcP} \quad \text{and} \quad \mathbf{C}^k = \text{diag} \left(\begin{matrix} 1 \\ \vdots \\ c_i^k \\ \vdots \\ \infty \end{matrix} \right)_{i=1, Nx+NcP} \quad (20)$$

378

379 The equilibrated matrix is defined at iteration $k + 1$ by

$$\tilde{\mathbf{Z}}^{k+1} = \mathbf{R}^k \cdot \tilde{\mathbf{Z}}^k \cdot \mathbf{C}^k \quad \dots \quad \dots \quad (21)$$

381 equal to 1 or after 50 iterations. Let \mathbf{R} and \mathbf{C}
 382 This procedure is repeated until all " r_i^k " and " c_i^k " are

383 403
 384 404
 385 405

$$\tilde{\mathbf{Z}} = \tilde{\mathbf{Z}} \cdot \mathbf{R} \cdot \mathbf{C}$$

$$= \tilde{\mathbf{Z}} \cdot \mathbf{R} \cdot \mathbf{C}$$

methods, such as Gaussian elimination [34] or LU decomposition [17, 40, 42]. In its actual form, the speciation code SPECY [48] uses unsymmetric multifrontal (UMF) [55] as the linear solver. To the best of our knowledge, no speciation code uses iterative methods to solve linear systems. This point is in accordance with the existing literature, which reports the use of iterative methods for solving large, sparse linear systems [20–22, 24, 26, 28, 29, 56, 57]. Nevertheless, actual developments in speciation codes involve the use of large chemical databases [39, 58, 59], leading to an increase in the size of the chemical systems. The use of iterative methods is also studied in this work.

We select some direct and iterative solvers according to the properties of the linear systems and the speciation computation methods currently in use (Table 1).

For the direct method, we select LU decomposition because it was originally used for speciation computations by Westall [40] and Westall et al. [42]. The UMF method has been implemented in the speciation code SPECY [48] in place of the LU approach [17]. After showing

that the Jacobian matrix is symmetric, we test the

DSYTRS subroutine from LAPACK [61], which is based on a UDU decomposition. Because the Jacobian matrix is

DPOTRS subroutine [61] based on the Cholesky method.

often positive definite, as shown in Table 3, we test the

ing preconditioning matrices and \mathbf{Z} the equilibrated matrix. Instead of solving the linear system (12), we solve

$$\mathbf{Z} \cdot \mathbf{x} = -\mathbf{Y} \quad (22)$$

where $\mathbf{x} = (\Delta\xi, \Delta Cp)$ and $\mathbf{Y} = \mathbf{R} \cdot \mathbf{Y}$. These procedures are coded using quadruple-precision reals. The linearsystem (22) is solved by LU decomposition coded with quadruple-precision real.

Even if the condition numbers of the Jacobian matrices (\mathbf{Z}) are very high ($10^{213.9}$ for the Fe-Cr mineral test case), the condition numbers of the equilibrated matrices (\mathbf{Z}) are much lower: the maximum condition number obtained after equilibration is $10^{13.4}$. According to Golub and van Loan [54], if the unit roundoff is approximately 10^{-d} and the condition number is approximately 10^q , then the Gaussian elimination gives a solution with approximately $d - q$ correct digits. Because we use quadruple precision, we obtain $d = 32$, leading to 32 - 14 = 18 correct digits. One can then assume that the reference solution is exact if we compare it to the solutions produced by the tested methods (computed using double-precision real).

2.6 Selected numerical methods for solving linearsystems

Studies on linear algebra [19, 23] present methods for solving linear systems as direct or iterative methods. Historically, speciation codes solved linear systems using direct

<p>332] Comput Geosci 333] of the Jacobi 334] [23, 335] Solv 336] react 337] tive 338] tran 339] spo 340] rt 341] und 342] er a 343] glo 344] bal 345] app 346] roa 347] ch. 348] Her 349] e, 350] we 351] test 352] QR 353] 432 354] dec 355] om 356] posi 357] tion 358] usin 359] g 360] the 361] DG 362] EL 363] S 364] rout 365] ine 366] [61] 367] . 368] 433 369] F 370] o 371] r 372] t 373] h 374] e 375] it 376] e 377] r 378] a 379] ti 380] v 381] e 382] m 383] e 384] t 385] h 386] o 387] d 388] s 389] , 390] w 391] e 392] t 393] e</p>	<p>434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452</p>	<p>the max of 8 ($N_x + N_cP$). The results obtained using the Jacobi and SOR methods are not detailed here. As previously reported [19], the Jacobi method is inefficient, leading to a very high failure ratio (close to 100 %) even for the easiest test cases. For the SOR (SOR) [23, 62] methods. Barrett et al. [21] proposed an algorithm to select an iterative solver depending on the matrix properties. GMRES was presented as the least selec- tive algorithm. We use a GMRES method developed by HSL [63]. If the matrix is symmetric, Barrett et al. [21] rec- ommend the use of conjugate gradient squared (CGS) or biconjugate gradient stabilized (BiCGStab) methods. CGS and BiCGStab subroutines have been developed by HSL. We test two additional methods devoted to symmetric matri- 444ces: SYMMBK [63] and an incomplete Cholesky (Inc. 445 CHOLESKY) factorization [63]. We use the same parameters for all iterative methods: a 447 maximum of 500 iterations and a stopping criterion of 10^{-8}. To determine the influence of the stopping criterion, we test the GMRES method using 50,000 maximum iterations and 10^{-12} as the stopping criterion, denoted by GMRES 10^{-12} in this study. A critical point of the GMRES algorithm is the size of the Hessenberg matrix. In this work, we set it to</p>	<p>453 454 455 456 457 458</p>
---	--	---	---

Table 1 List of the selected solvers

Name	Source	Method	Matrix properties
Direct			
LU	[58]	LU decomposition	–
DGETRS	[59]	LU decomposition	–
UMF	[53]	Direct multifrontal	–
DSYTRS	[59]	UDU-factored symmetric matrix	Symmetric
DPOTRS	[59]	Cholesky $A = U^T \times U$	Definite positive
DGELS	[59]	QR decomposition	–
LU QUAD	[58]	LU decomposition quadruple precision	–
Iterative			
SYMMBK	[61]	Iterative SYMMBK HLS_MI02	Symmetric
Inc. CHOLESKY	[61]	Incomplete Cholesky HSL_MI28	Symmetric
CGS	[61]	Conjugate gradient squared HLS_MI23	–
BiCGStab	[61]	Biconjugate gradient squared stabilized HLS_MI26	–
GMRES	[61]	Flexible GMRES HLS_MI15	–
Gauss-Seidel	[58]	Gauss-Seidel method	–
Preconditioned			
LU Equil	[51–58]	LU and matrix equilibration	–
DGESVX	[59]	LU and optional preconditioning	–
GMRES Equil	[51–61]	GMRES and matrix equilibration	–
GMRES 1.d-15	[61]	GMRES convergence criteria 1.d-15	–

459 method [23, 26, 56, 62], the over-relaxation parameter is the
 460 key factor. Unfortunately, we did not find any efficient relations-
 461 hips to define it. For the same chemical system, the best
 462 value varies from 0.097 to 1.91 without apparent order.

463 We do not extensively test the possibility of using a pre-
 464 conditioner. As stated by Barrett et al. [21]: “Since using a
 465 preconditioner in an iterative method incurs some extra cost,
 466 both initially for the setup, and per iteration for applying
 467 it, there is a trade-off between the cost of constructing and
 468 applying the preconditioner, and the gain in convergence
 469 speed”. In our case, the matrices are very small, leading us

470 to suppose that this trade-off would not be advantageous.
 471 Nevertheless, an easy way to test preconditioners is pro-
 472 posed by the LAPACK routine DGESVX, which performs
 473 LU decomposition and matrix equilibration depending on
 474 the estimated condition number. We implement matrix equi-
 475 libration according to Knight et al. [53] to obtain a reference
 476 solution. We test this preconditioning technique associated
 477 with LU decomposition and the GMRES method, denoted
 478 by LU Equil and GMRES Equil in this study. The maximum
 479 iterations allowed for the equilibration procedure is fixed to
 480 5, according to the recommendations of Knight et al.

Table 2 Structure of the Jacobian matrix

Z	$\frac{\partial}{\partial \xi_k}$	$\frac{\partial}{\partial [C_{p_{k-N_X}}]}$
$\frac{\partial Y_j}{\partial}$	$\sum_{i=1}^{N_C} a_{i,j} \cdot a_{i,k} \cdot [C_i^n]$	$ap_{k-N_X,j}$
$\frac{\partial Y_j}{\partial}$	$\sum_{i=1}^{N_C} a_{i,j} \cdot a_{i,k} \cdot [C_i^n]$	0
	$ap_{k,j-N_X}$	

Table 3 Properties of the 10 chemical test cases ranked by increasing the maximal condition number

	Nx	Nc	NcP	Z size	cond(Z) min	cond(Z) max	cond(Z) max after 20 equilibration	%Z diagonal dominant	%Z positive definite
Gallic acid	3	17	0	3	10 ^{0.61}	10 ^{12.6}	10 ^{0.95}	18.4	100
Valocchi	5	7	0	5	10 ^{0.49}	10 ^{15.3}	10 ^{0.65}	67.7	100
Pyrite	4	40	0	4	10 ^{4.06}	10 ^{24.9}	10 ^{0.95}	0.00	100
MoMaS easy	5	12	0	5	10 ^{3.44}	10 ^{37.7}	10 ^{1.05}	0.00	71.1
Morel-Morgan	52	781	0	52	10 ^{43.4}	10 ^{60.7}	10 ^{1.13}	0.00	35.9
MoMaS medium	5	14	0	5	10 ^{5.88}	10 ^{103.9}	10 ^{0.95}	0.00	78.8
Fe-Cr	7	39	0	7	10 ^{9.46}	10 ^{113.6}	10 ^{1.05}	0.00	68.9
Pyrite mineral	4	43	3	7	10 ^{1.71}	10 ^{33.1}	10 ^{3.19}	0.00	0.00
MoMaS hard	6	15	2	8	10 ^{5.45}	10 ^{123.9}	10 ^{3.02}	0.00	0.00
Fe-Cr mineral	7	43	3	10	10 ^{8.67}	10 ^{213.9}	10 ^{13.4}	0.00	0.00

481 Finally, we test an LU decomposition method compiled
 482 as quadruple precision, denoted by LU QUAD. The source
 483 of this method is the LU double-precision real of numerical
 484 recipes [60], and we adapt it to quadruple precision.
 485 Because the usual computations are performed using double
 486 precision, the quadruple precision ($d\xi_{\text{QUAD}}$) should be
 487 translated in double-precision real. To avoid overflow, we
 488 rescale $d\xi_{\text{QUAD}}$ to ensure its validity. If huge (1.d0) is the
 489 highest double-precision real represented by the machine,
 490 we rescale $d\xi_{\text{QUAD}}$ to obtain the double-precision solution
 491 $d\xi_{\text{LU QUAD}}$:

$$d\xi_{\text{LU QUAD}} = \text{huge}(1.d0) \cdot d\xi_{\text{QUAD}} \quad (23)$$

492 In this way, we conserve the direction of the Newton step,
 493 even if its norm is changed.

494 3 Results and discussion

495 3.1 Properties of the Jacobian matrices

496 As defined by Eq. 11, the Jacobian matrix has several
 497 properties:

- 498 (i) The matrix is block-structured, as presented in Table
 499 2. A four-block structure is present if precipitation
 500 occurs.
- 501 (ii) The matrix is symmetric, as shown in Table 2.
- 502 (iii) In the case of no precipitation, all the diagonal terms
 503 of the matrix are strictly positive because they are the
 504 sum of $\sigma_{i,j}^2[C_i]$. It is then possible for the matrix to
 505 be diagonal dominant. We examine this possibility for
 506 the selected test case. Table 3 shows the ratio of diagonal
 507 dominant Jacobian matrices for all the chemical
 508 tests performed according to the previously defined

test procedure. Some matrices in the gallic acid and 509
 Valocchi cases are diagonal dominant, but none of the 510
 matrices from the other cases are diagonal dominant. 511
 By plotting the ratio of diagonal dominant matrices 512
 depending on the condition number (see Appendix 2 513
 (B-1)), it appears that only matrices with very low 514
 condition numbers can be diagonal dominant. 515
 (iv) Because the Jacobian matrix is real, symmetric, 516
 and sometimes diagonal dominant, the question of 517
 whether it is positive definite may be posed. In the 518
 case of no precipitation, Eq. 11 can be written in 519
 matrix form, leading to Eq. 24 520

$$Z = A^T \cdot \text{diag}(C) \cdot A \quad (24)$$

Because the concentrations are positive, the Jacobian matrix 521
 is analytically positive definite. Nevertheless, this may not 522
 be true numerically. We are not able to propose a gen- 523
 eral framework, but we can compute the eigenvalues of the 524
 Jacobian matrix and test whether they are positive for all 525
 cases. Table 3 shows that for the gallic acid, Valoc- 526
 chi, pyrite, and Morel-Morgan test cases, all the Jacobian 527
 matrices are positive definite. For the MoMaS easy, MoMaS 528
 medium, and Fe-Cr test cases, a large proportion (66.4 to 529
 74.1 %) of the Jacobian matrices are positive definite. For 530
 cases including minerals (pyrite mineral, MoMaS hard, and 531
 Fe-Cr), essentially none of the matrices are positive defi- 532
 nite (only 0.1 % for the MoMaS hard test). Plotting the ratio 533
 of positive definite matrices as a function of the condition 534
 number (see Appendix 2 (B-2)) shows that the chemical 535
 conditions are more important than the condition number 536
 when determining whether the Jacobian matrix is diagonal 537
 dominant. 538

- (v) According to the test procedure presented previously, 539
 we plot, on the same graph, the logarithm of the norm 540

541 of $\|Y\|$ and the logarithm of the condition number of the
 542 matrix Z (Fig. 1). There is a strong linear relationship
 543 between these parameters. Moreover, the linear
 544 relationship does not depend on the chemical test, only
 545 on the existence of minerals. According to our results,
 546 the conditioning of the Z matrix can be evaluated using
 547 the following empirical formulas:

$$\begin{aligned} \text{cond}(Z)_{\text{no mineral}} &= 10^{5.30 \pm 0.03} \times \|Y\|^{0.9374 \pm 0.0008} \\ \text{cond}(Z)_{\text{mineral}} &= 10^{-3.23 \pm 0.08} \times \|Y\|^{1.706 \pm 0.002} \end{aligned} \quad (25)$$

548 The value and uncertainties are obtained through the least
 549 squares method over all $\text{cond}(Z)$ and $\|Y\|$. In this way, we
 550 propose an estimation of $\text{cond}(Z)$ with no computation time
 551 cost because the objective function is evaluated during the
 552 Newton-Raphson procedure. As shown in Fig. 1, $\text{cond}(Z)$
 553 and $\|Y\|$ are strongly correlated for large condition numbers,
 554 and the results are noisier if $\text{cond}(Z)$ and $\|Y\|$ are

small. The evolution of this relation for low $\|Y\|$ can be seen
 Appendix 7 (G-11). Therefore, Eq. 25 should not be used
 for $\|Y\|$ s than 10^{10} .

Several of these properties are obtained using the logarithm of the component activities as the primary unknown in Eq. 8. The historical approach [34] uses the component concentrations as the primary variable and leads to a less interesting Jacobian matrix. Even if the structure presented in Table 2 exists, the matrix is not symmetric. Moreover, the matrix is worse conditioned (condition number from $10^{11.2}$ to $10^{49.4}$ rather than $10^{4.06}$ to $10^{24.9}$ for the pyrite case). Finally, no specific relation exists between $\text{cond}(Z)$ and $\|Y\|$ for the historical formulation.

As an example, we show one linear system from the Fe-Cr mineral test, corresponding to a condition number of 10^{187} . One can observe the structure of the matrix and the specificity of the linear system (26).

$$\begin{pmatrix} 1.15 \cdot 10^{94} & 9.09 \cdot 10^{93} & -5.04 \cdot 10^{-13} & -11.7 & 3.03 \cdot 10^{93} & 0 & 1.10 \cdot 10^{87} & 0 & 5 & -1 \\ & 5.45 \cdot 10^{93} & 1.00 \cdot 10^{-10} & 2.374 \cdot 10^{10} & 1.82 \cdot 10^{93} & 0 & 4.11 \cdot 10^{86} & 0 & 0 & 0 \\ & & & 2.91 & 8.74 & 0 & 1.28 \cdot 10^{-6} & 1 & 0 & 0.25 \\ & & & & 6.06 \cdot 10^{92} & 2.285 \cdot 10^{-22} & 1.371 \cdot 10^{86} & 0 & 0 & 0.75 \\ & & & & & & 1.37 \cdot 10^{86} & 0 & 0 & 0 \\ & & & & & & & 0 & 0 & 0 \\ & & & & & & & & 0 & 0 \\ & & & & & & & & & 0 \end{pmatrix} (d\xi) = \begin{pmatrix} -3.03 \cdot 10^{93} \\ -1.82 \cdot 10^{94} \\ -1.05 \cdot 10^{-13} \\ 8.99 \cdot 10^{-3} \\ -6.06 \cdot 10^{94} \\ 2.25 \cdot 10^{-2} \\ -1.37 \cdot 10^{86} \\ -27.6 \\ 180 \\ 3.84 \end{pmatrix} \quad (26)$$

572 **3.2 Robustness of the methods**

573 Figure 2 presents the failure ratio for each method and each
 574 test case. The presence of minerals prevents the DPOTRS,
 575 Inc. CHOLESKY, and Gauss-Seidel methods from solving
 576 the system. If there are minerals present in the chemical
 577 system, a zero-value block appears in the Jacobian matrix,
 578 as shown in Table 2 and Eq. 26. This block makes the
 579 Inc. CHOLESKY factorization unappropriated. Because the
 580 Gauss-Seidel method requires division by each diagonal
 581 term, this zero-value block makes the method unadapted.
 582 The failure of the DPOTRS routine is explained by the
 583 properties of the Jacobian matrix. As shown in Table 3,
 584 there is no positive definite matrix in the presence of minerals.
 585 In the case of the DPOTRS, Inc. CHOLESKY, and
 586 Gauss-Seidel methods, the term failure is ambiguous. These
 587 methods are *expected to fail* and should not be used on
 588 systems with minerals. If there are no minerals, some matrices
 589 are not positive definite in the MoMaS easy, MoMaS
 590 medium, Morel-Morgan, and Fe-Cr tests. This explains the

Some other methods (DGETRS, DSYTRS, DGELS, and DGESVX) present a substantial failure ratio, mainly for high condition number tests (MoMaS easy and Fe-Cr mineral). UMF, SYMMBK, and CGS are robust for the Fe-Cr mineral test but present significant failure ratios for lower-conditioned tests, such as MoMaS easy or pyrite mineral. Some methods adapted to symmetric matrices (DSYTRS and SYMMBK) are included in this class of weak methods.

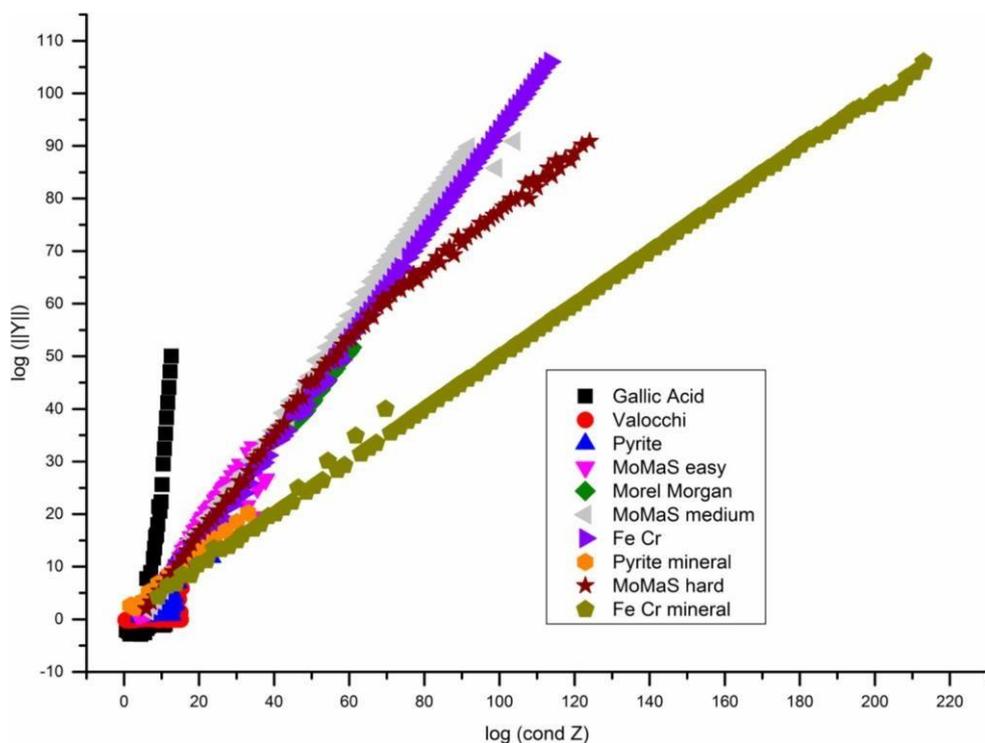
The BiCGStab method has a very low failure ratio and fails only in the two difficult tests (MoMaS easy and Fe-Cr mineral). GMRES is the only successful iterative method.

Figure 2 shows that some methods are successful for all the test cases. The most successful direct method is LU, while the most successful iterative methods are GMRES and GMRES 10^{-12} . The quadruple-precision method LU QUAD is also successful, which is expected because the double-precision LU method is also successful. The use of an equilibration method as a preconditioner makes LU Equil and GMRES Equil successful.

As stated previously, we focus on the capacity of a

591 ^{Comput Geosci} failure of the DPOTRS routine. method to produce a solution independent of its accuracy. ^{Comput Geosci} 612

Fig. 1 Relationship between the condition number of Z and the norm of the objective function plotted on a log-log scale



613 For some advanced methods (e.g. LAPACK methods), a
 614 a posteriori estimation of the residual and estimation of the
 615 condition numbers are performed. If the solution is not suf-
 616 ficiently accurate, no solution is given, leading to a higher

failure ratio than for the more rustic methods (LU or Gauss-
 Seidel). Because the key point of this work—the resolution
 of a linear system—is included in the iterative Newton pro-
 cedure, it is preferable to obtain an inaccurate solution (so

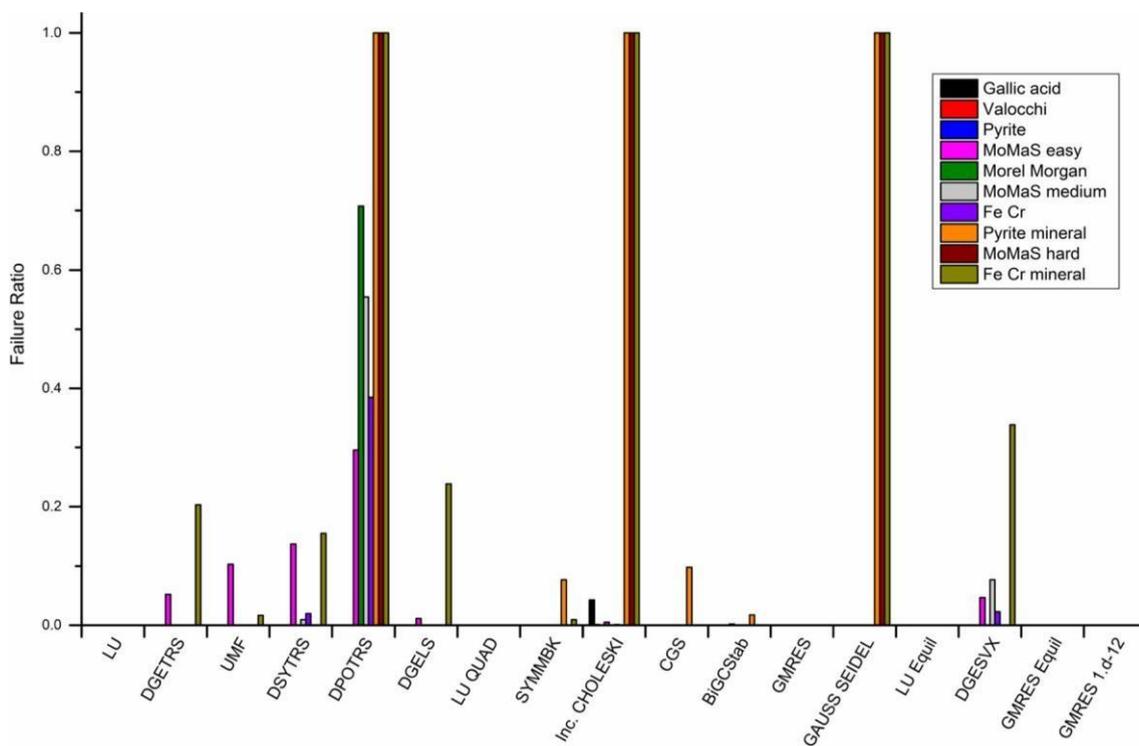


Fig. 2 Mean of the failure ratio for each method and each test case

621 the iterative procedure can be continued) than no solution
 622 (the iterative procedure will be aborted).

623 Appendix 3 presents the evolution of the failure ratio for
 624 each test case and each method depending on the condition
 625 number.

626 For the direct methods (Appendix 3 (C-1 to C-5)), for
 627 small condition numbers corresponding to the test cases gal-
 628 lic acid, Valocchi, and pyrite, no failure occurs. As the con-
 629 dition number increases, the failure ratio also increases for
 630 some methods. MoMaS easy (Appendix 3 (C-4)), MoMaS
 631 medium (Appendix 3 (C-6)), and Fe-Cr (Appendix 3 (C-
 632 7)) show that for condition numbers greater than 10^{20} , the
 633 failure ratio increases greatly for some of the methods.
 634 These methods are DOPTRS and DSYTRS for MoMaS
 635 medium and Fe-Cr. DGETRS, UMF, DSYTRS, DOPTRS,
 636 and DGELS present some failure for condition numbers
 637 greater than 10^{15} for the MoMaS easy case. In the presence
 638 of minerals (Appendix 3 (C-8 and C-9)), for low condition
 639 numbers (the pyrite mineral case), the methods are either
 640 successful (UMF, LU, DSYTRS, DGETRS) or completely
 641 unsuccessful (DPOTRS). For very high condition numbers
 642 (Fe-Cr mineral case), the success of the method does not
 643 depend on the condition number. We suppose that the con-
 644 dition numbers (see Table 3) are too high to exhibit any
 645 ordering.

For other iterative methods, the success does not depend
 on the condition number but on the nature of the matrix
 and the presence (Appendix 3 (C-18 to C-20)) or absence
 (Appendix 3 (C-11 to C-17)) of minerals.

3.3 Accuracy of the methods

The accuracy of the methods is evaluated in two ways: (i)
 the relative error on the norm (18) and (ii) the angle between
 the reference and the calculated solution (19).

(i) By plotting the mean of the logs of the relative error
 the norm of each test case (Fig. 3), some general
 tendencies are identified. The relative residual tends to
 increase with the condition number of the system. For
 direct methods and small condition numbers, the rel-
 ative residual is small (10^{-10} to 10^{-3}) for the gallic
 acid, Valocchi, and pyrite test cases. For the itera-
 tive methods, the relative residual corresponding to an
 accurate resolution for tests with small condition num-
 bers is approximately 10^{-4} . This value corresponds to
 the value of the convergence criteria of the iterative
 methods. Iterative methods are more sensitive to the
 condition number than direct methods. Only the Val-
 occhi test case is accurately solved by almost all the

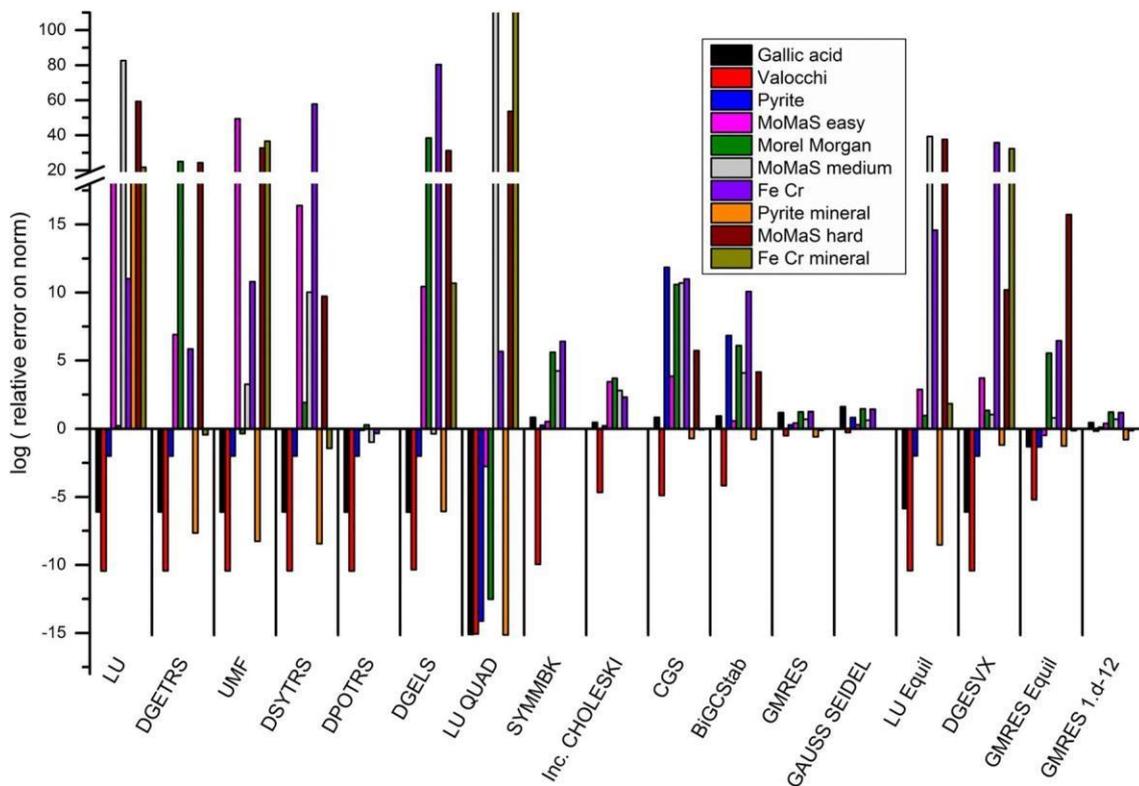


Fig. 3 Mean of the logs of the relative error on norm for each method and each test case

668 iterative methods, whereas the first three tests are accu- 692
 669 rately solved by all the direct methods. Even in the case 693
 670 of successful resolution (CGS and BiCGStab methods), 694
 671 the relative errors on the norm are high for intermediate 695
 672 cases (pyrite, MoMaS easy, and Morel-Morgan). Nev- 696
 673 ertheless, the results are better for the iterative meth- 697
 674 ods than for the direct methods for the difficult tests 698
 675 (MoMaS easy, MoMaS medium, MoMaS hard, Fe- 699
 676 Cr mineral). The GMRES and Gauss-Seidel methods 700
 677 have mostly constant mean relative error on the norm, 701
 678 with the same accuracy for all test cases. GMRES and 702
 679 Gauss-Seidel are less efficient than the other methods 703
 680 for the easy tests, but more ill-conditioned tests are 704
 681 better solved by these two methods. 705

682 The condition numbers are so high that even LU QUAD 706
 683 cannot provide accurate resolution. For the MoMaS medium 707
 684 and Fe-Cr mineral tests, many of the solutions calculated 708
 685 by the LU QUAD method are rescaled using Eq. 23, leading 709
 686 to excessively high relative error on the norm.

687 Comparison of the relative error on the norm given by 710
 688 the non-preconditioned (LU, DGETRS, and GMRES) and 711
 689 preconditioned (LU Equil, DGESVX, and GMRES Equil) 712
 690 methods shows that the preconditioned methods lead to 713
 691 lower relative error than the non-preconditioned methods 714
 715
 716

for the direct methods, but the result is more case-dependent 692
 for GMRES. The use of preconditioning usually leads to 693
 lower relative error on the norm, except for the Morel- 694
 Morgan, Fe-Cr, and MoMaS hard cases. 695

Increasing the maximum number of iterations and reduc- 696
 ing the convergence criteria of GMRES leads to less relative 697
 error on the norm, but this reduction is not significant. 698

Nevertheless, the global means of the logs of relative 699
 errors on the norm hide the influence of the increasing 700
 condition number. Appendix 4 presents the evolution of 701
 the relative error on the norm for each test case and each 702
 method depending on the condition number. The theoret- 703
 ical behaviour is verified for the direct methods and for all 704
 the test cases (except for the Valocchi one, Appendix 4 (D- 705
 2)). The relative error on the norm increases regularly with 706
 the condition number. It is close to 10^{-16} when the con- 707
 dition number is close to 1 and increases to 1 when the 708
 condition number is close to 10^{16} , in accordance with the 709
 computation theory presented by Golub and van Loan [54]. 710
 For condition numbers greater than 10^{16} , the evolution of 711
 the relative error on the norm with the condition number is 712
 much noisier. The use of the quadruple-precision LU QUAD 713
 method leads to an accurate resolution of a large portion 714
 of the tested systems. As expected by computation theory, 715
 all the systems with condition numbers less than 10^{32} are 716

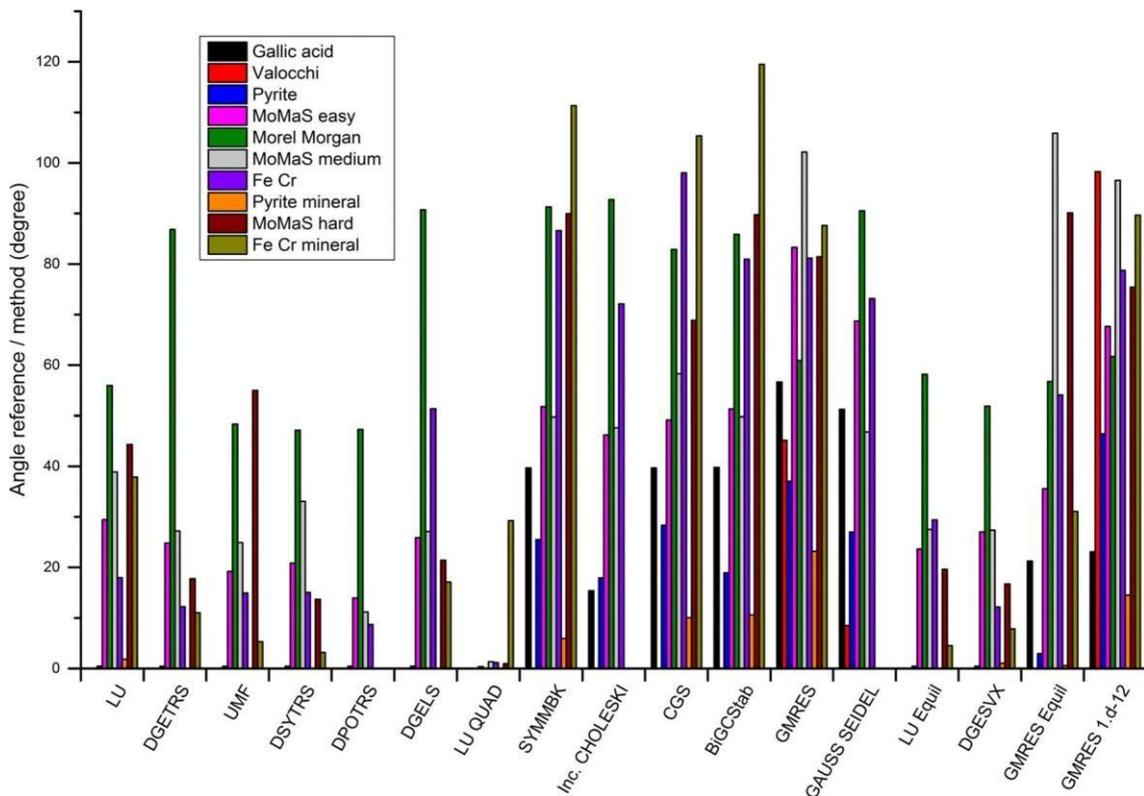
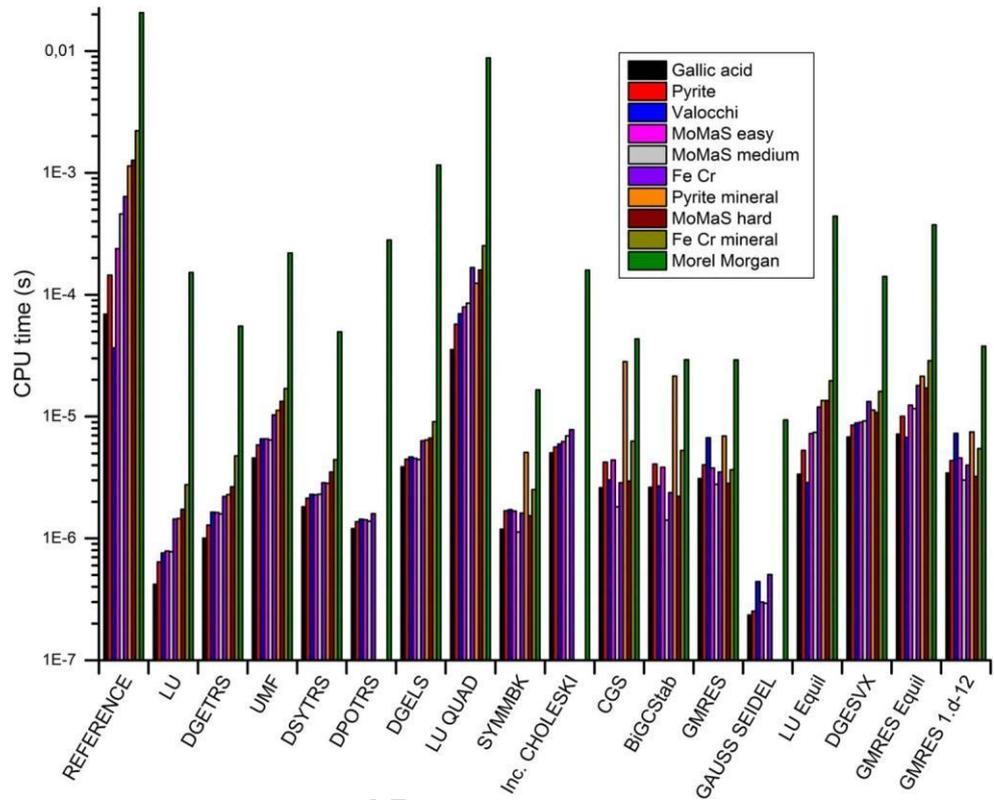


Fig. 4 Mean of the angles between reference and computed solution for each method and each test case

Q7

Fig. 5 CPU time (s) for each method and each test case



717 solved with a relative error on the norm of approximately
 718 10^{-15} . In some cases (MoMaS medium (Appendix 4 (D-
 719 6)), Fe-Cr (Appendix 4 (D-7)), MoMaS hard (Appendix 4
 720 (D-9))), LU QUAD produces an increasing relative error
 721 with increasing condition number (if higher than 10^{32}) but
 722 not systematically. LU QUAD produces a very low rela-
 723 tive error on the norm even if the condition number is very
 724 high (Appendix 4 (D-9)). This behaviour can be explained
 725 by the fact that the LU QUAD method and/or the reference
 726 method is unable to exactly solve such ill-conditioned sys-
 727 tems. LU QUAD produces a very high relative error on the
 728 norm, one point with 10^{290} error for the MoMaS medium
 729 (Appendix 4 (D-6)), and all the values at condition numbers
 730 greater than 10^{90} for the Fe-Cr mineral (Appendix 4 (D-10))
 731 test case. These points correspond to the rescaling of the
 732 computed quadruple-precision solution to maintain it on the
 733 double-precision scale (using Eq. (23)).

734 Iterative methods present similar behaviour to direct
 735 methods, giving very low relative error on the norm
 736 (between 10^{-15} and 10^{-8}) when the condition number is
 737 less than a critical value. This critical value depends on the

pyrite (Appendix 4 (D-13)), the MoMaS medium (Appendix 744
 4 (D-16)), and MoMaS hard (Appendix 4 (D-19)) tests. 745
 Using low convergence criteria (GMRES 1.d-12) leads to 746
 lower relative error on the norm for low condition numbers 747
 (Appendix 4 (D-21 to D-23, D-26 to D-29)), but no sig- 748
 nificant improvements are obtained if the condition number 749
 increases, as shown in Appendix 4 (D-24 to D-30). 750

Using preconditioning methods reduces the relative error 751
 on the norm for intermediate condition numbers. No gain is 752
 obtained for low condition numbers (Appendix 4 (D-21 and 753
 D-22)), but the errors given by LU Equil, DGESVX, and 754
 GMRES Equil are less than the LU and GMRES errors for 755
 higher condition numbers (Appendix 4 (D-24 to D-26)). For 756
 very high condition number tests (Appendix 4 (D-27, D-29, 757
 and D-30)), the errors given by the preconditioned methods 758
 are equivalent to the errors given by the non-preconditioned 759
 methods. 760

- (ii) By plotting the angle between the reference solu- 761
 tion and the calculated solution, we can compare the 762
 methods according to the computed direction (Fig. 4). 763

Because the resolution of the linear system (13) repre- 764 Q8

method and the test case. It can be set to 10^8 for SYMMBKCGS,
 BiCGStab, and GMRES for the gallic acid (Appendix4 (D-11)) and
 MoMaS easy (Appendix 4 (D-14)) cases. It can be set to 10^{12} or
 10^{15} for Inc. CHOLESKY for the gallic acid and MoMaS easy cases
 and for SYMMBK, Inc. CHOLESKY, CGS, BiCGStab, and

738 743
 739
 740
 741
 742

GM ^{Comput Geosci}
RES ~~sents one step~~ in the iterative Newton procedure.
for this ⁷⁶⁵ information is much more important than
the the norm ⁷⁶⁶of the step. A wrong norm can be
corrected using ⁷⁶⁷ line search methods [64],
whereas modifying a wrong ⁷⁶⁸ direction leads to
additional iterations. Small condition ⁷⁶⁹ number
tests (gallic acid, Valocchi, pyrite, and pyrite ⁷⁷⁰

Comput Geosci

771 mineral) are solved using direct methods with the right
 772 direction. If the condition number increases, the direc-
 773 tions given by the direct methods become inaccurate,
 774 but the condition number is not the only govern- ing
 775 parameter. Morel-Morgan leads to worse direction than
 776 MoMaS medium and Fe-Cr, and MoMaS hard leads to
 777 a higher angle than the Fe-Cr mineral test. Iter-ative
 778 methods result in a worse direction than direct
 779 methods, and only the Valocchi test case is solved with
 780 an accurate direction by all the iterative meth- ods.
 781 Imposing lower convergence criteria (10^{-12}) on
 782 GMRES leads to a worse direction than using the usual
 783 criteria (10^{-8}). Using preconditioning methods leads to
 784 a better direction when associated with a direct method
 785 (LU Equil and DGESVX), but the conclusion is less
 786 clear for the iterative GMRES Equil method.
 787 Depending on the test case, the direction can be worse
 788 (Valocchi, MoMaS easy, MoMaS medium) or better
 789 (gallic acid, pyrite, MoMaS hard, Fe-Cr mineral)

790 The influence of the condition number on the angle (see
 791 Appendix 5) indicates that the direction is correct for direct
 792 methods when the condition number is less than 10^{15} . For

793 iterative methods, the limit to obtain an accurate direction is
 794 a condition number less than 10^8 , excepted for the Gauss-
 795 Seidel method, which produces wrong directions for low
 796 condition numbers. If the condition number increases, the
 797 behaviour of the direction becomes noisy. Since the rela-
 798 tive error on the norm increases regularly until the condition
 799 number reaches the limit of 10^8 or 10^{15} , the angle is accu-
 800 rately defined until this condition number limit is reached.
 801 Using preconditioned methods leads to a better direction for
 802 the LU Equil and the GMRES Equil methods when the con-
 803 dition number is higher than 10^{15} for some cases (Appendix

804 5 (E-21, E-23 to E-25, and E-30)) but to a worse direction
 805 for other cases (Appendix 5 (E-26 and E-29)).

806 We present two successful direct methods, LU and LU
 807 QUAD; one iterative method, GMRES (both tested ver-
 808 sions, GMRES and GMRES 10^{-12}); and two precondi-
 809 tioned methods, LU Equil and the GMRES Equil. By compar-
 810 ing the relative error on the norm (Appendix 4 (D-21 to
 811 D-30)), the successful methods can be ranked from the low-
 812 est to highest error: LU QUAD, GMRES 10^{-12} , GMRES
 813 Equil, LU Equil, and LU. Ranking these methods according
 814 to the angle between the reference and computed solution is

more complicated. For all the tests cases (Appendix 5 (E- 815
 21 to E-25, E-27, E-28, and E-30)), LU QUAD gives the 816
 best direction, followed by LU Equil, LU, GMRES Equil, 817 and
 GMRES 10^{-12} . The MoMaS medium (Appendix 5 (E- 818 26)) and
 MoMaS hard (Appendix 5 (E-29)) test cases lead 819 to the same
 conclusion, except GMRES Equil which gives 820 the worst
 direction. 821

3.4 Efficiency of the methods 822

The speed of the methods is studied by recording the com- 823
 putation time for each test case and plotting the mean CPU 824
 time for each test case and each method (see Fig. 5). As 825
 expected, the computation times are very short (less than 1 826
 ms) because the systems to solve are small. 827

Figure 5 shows the influence of the system size. For all 828
 methods, the computation time increases with the number 829
 of unknowns. The results show that the iterative methods 830
 are less sensitive to the system size than the direct meth- 831
 ods. For the iterative methods, the number of iterations is 832
 important and depends on the first guess and other factors. 833
 The slowest method is LU QUAD, for which a large 834
 amount of computation time is devoted to the translation of

double-precision real to quadr. ple-precision real and back. 835
 Figure 5 also shows the computing time required to obtain 836
 the reference solution, which requires more time. 837

The UMF method is the slowest double-precision direct 838
 method, but its multifrontal block strategy becomes interest- 839
 ing for large system. The resolution of the Morel-Morgan 840
 test requires 33 times more CPU time than the resolution of 841

the MoMaS easy test for the UMF method, whereas it takes 842
 190 times more time for the LU method. 843
 844

Among the iterative methods, the fastest is the Gauss- 845
 Seidel method and the slowest is the the CHOLESKY 846

method. The two most robust iterative methods, BiCGStab 847 and
 GMRES, are rapid, sometimes more so than the direct 848 robust
 methods, LU and UMF, especially for large systems 849 (Morel-
 Morgan test case). GMRES is less case-dependent 850 than
 BiCGStab, leading to similar computing time, regard-

10 ⁴ < cond(Z) < 10 ¹⁰	LU Nx	NcP < 10	GMRES Nx	NcP 10
10 ⁴ > cond(Z)	LU			

Table 4 Algorithm for equilibrium computation

cond(Z)	Inversion method
>10 ³⁰	LU QUAD Equil
10 ³⁰ ≥ cond(Z) > 10 ¹⁴	LU QUAD
10 ¹⁴ ≥ cond(Z)	+
+	+
+	+
+	+

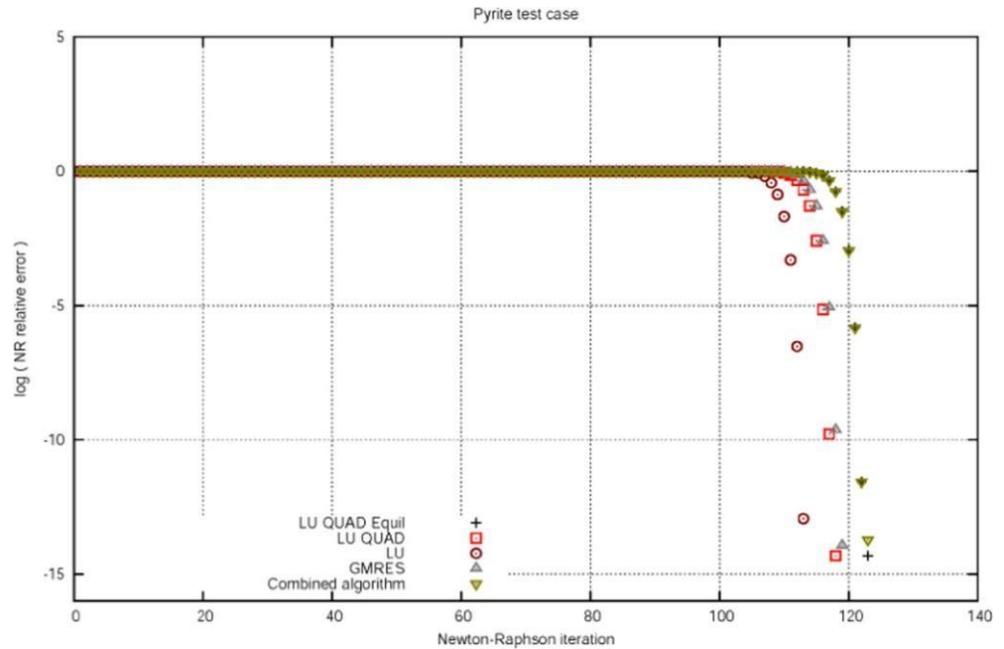
less
of
the
test
cas
e. 85
2

A
s
e
x
p
e
c
t
e
d
,
i
n
t
r
o
d
u
c
i
n
g
p
r
e
c
o
n
d
it
i
o
n
i
n
g
t
e
c
h
n
i
q
u
e
s
8
5
3

(LU Equil, DGESVX, and GMRES Equil) or decreasing 854
the convergence criteria for an iterative method (GMRES 855
 10^{-12}) leads to increased computing time. The computing 856
time for preconditioning does not depend only on the sys- 857
tem's size: the Valocchi, MoMaS easy, and MoMaS medium 858
test cases (system size of 5×5) are solved with the same 859
computing time for all the direct methods, but their resolu- 860
tion when using LU QUAD Equil, LU Equil, and GMRES 861
Equil is faster. 862

Appendix 6 shows the computation time (log scale) for 863 Q9
each test case and each method depending on the condition 864
number. Appendix 6 (F-1 to F-10) shows that, as expected, 865

Fig. 6 Evolution of $NR_{\text{relative error}}$ as a function of the Newton-Raphson iteration for the pyrite test case



866 the computation time of the direct methods does not depend
 867 on the condition number of the system. The LU method is
 868 usually 10 times faster than the UMF method, except for the
 869 Morel-Morgan test case, in which LU is only 1.5 times faster.

870 In Appendix 6 (F-11 to F-20), the general tendency for
 871 the iterative methods is to require the same computation
 872 time, independent of the condition number. The oscillations
 873 presented by the curves seem to be not related to the
 874 condition number. For the test case without minerals,
 875 the Gauss-Seidel method is efficient. The two most robust
 876 methods, BiCGStab and GMRES, are often the third and
 877 fourth fastest methods (Gauss-Seidel and SYMMBK are the
 878 fastest).

879 **4 Proposal of a new algorithm**

880 Based on our results, we propose an algorithm to opti-
 881 mize the resolution of a chemical system using a Newton-
 882 Raphson-like method.

883 Examining the failure ratio results, seven methods are
 884 eligible: LU and LU QUAD as direct methods, GMRES
 885 and Gauss-Seidel (if no minerals) as iterative methods, LU
 886 Equil and GMRES Equil as preconditioned methods, and
 887 the reference method (LU QUAD Equil).

888 Because these methods are included in a Newton mini-
 889 mization procedure, the most important accuracy criterion
 890 is the direction of the minimization, i.e. the angle between
 891 the reference and the calculated solution. The behaviour of
 892 this direction is strongly correlated with the condition num-
 893 ber of the system and is correct if the condition number
 894 is less than the critical value and wrong if the condition

number is greater than the critical value (see Appendix 5). 895
 The critical condition number is 10^8 for GMRES, 10^{16} for 896
 the double-precision direct methods, 10^{32} for LU QUAD, 897
 and case-dependent for preconditioned methods (10^{20} to 898
 10^{60}). Gauss-Seidel leads to wrong directions for very low 899
 condition numbers (Appendix 5 (E-11 and E-12)). 900

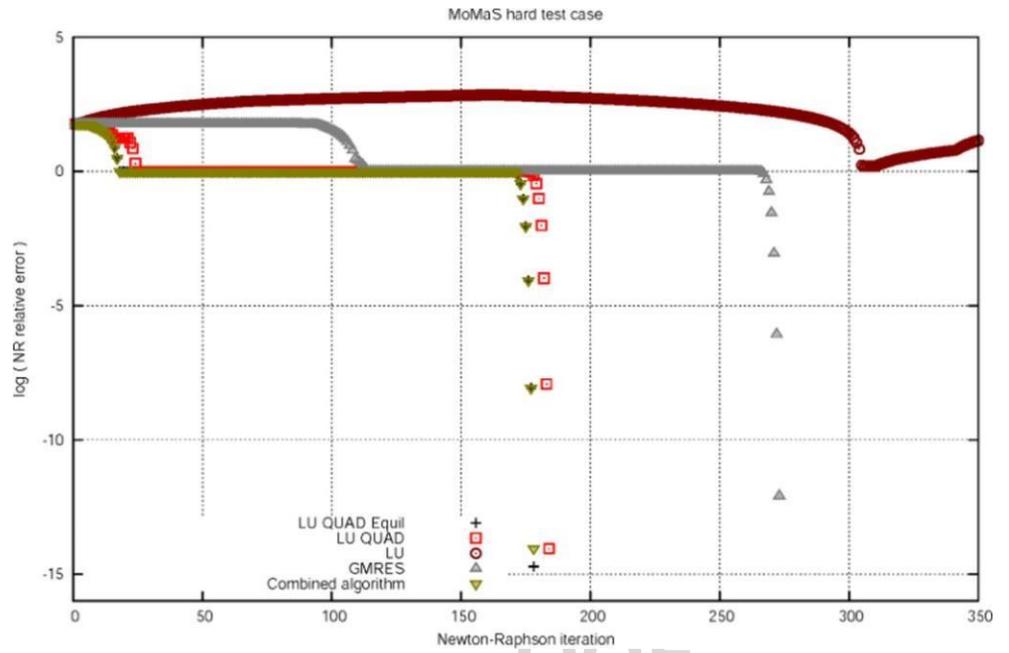
In terms of efficiency, the most rapid method is Gauss- 901
 Seidel when it is available. The second most efficient 902
 method is LU for small systems (less than 10×10) or 903
 GMRES for larger systems (more than 10×10), and the 904
 slowest method is LU QUAD. For small systems (less than 905
 5×5), LU Equil is as fast as GMRES but becomes slower as 906
 the system size increases. 907

We recommend using LU, LU QUAD, GMRES, and the 908
 reference method LU QUAD Equil. Gauss-Seidel should be 909
 rejected because of its wrong direction, and equilibration 910
 does not sufficiently improve the behaviour of double- 911
 precision routines. 912

Using Eq. 25, it is possible to estimate the condition 913
 number of the system without additional computation. This 914
 estimation enables the selection of the best-adapted method 915
 depending on the system size and condition number. 916

The goal is to use the most robust method (LU QUAD 917
 with preconditioning) for high condition number systems 918
 (more than 10^{32}) in the first Newton-Raphson iterations. 919
 When the condition number is sufficiently decreased, the 920
 preconditioning becomes useless and LU QUAD can be 921
 used until the condition number is less than 10^{16} . Then, a 922
 faster method is used to obtain a coarse approximation of 923
 the solution, LU for small systems and GMRES for large 924
 systems (more than 10×10). To find the exact solution, the 925
 LU direct method is used. 926

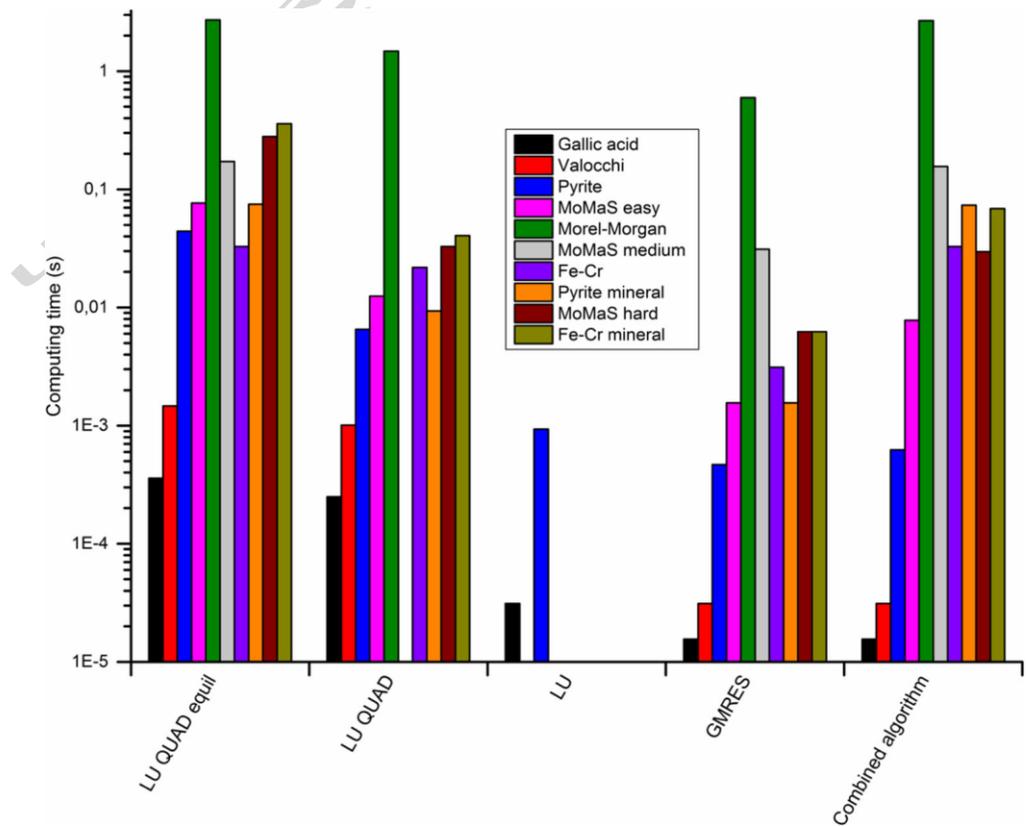
Fig. 7 Evolution of $NR_{relativeerror}$ as a function of the Newton-Raphson iteration for the MoMaS hard case



927 We propose the algorithm presented in Table 4 and
 928 compare it with several inversion methods in a Newton-
 929 Raphson algorithm. The 10 chemical test cases are solved
 930 using the combined algorithm or one of the selected meth-
 931 ods: LU QUAD Equil (used as the reference solution), LU
 932 QUAD, LU, and GMRES. Appendix 7 shows the evolution...

of the $NR_{relativeerror}$ (7) as a function of the Newton-Raphson
 iterations
 Figure 6 shows that all the methods are equivalent for
 easy test cases (see Appendix 7 (G-1 to G-3)). Nevertheless,
 the use of LU inversion leads to non-convergence, even if
 the test is easy, as observed for the Valocchi test (Appendix

Fig. 8 Computation time (s) as a function of test case and algorithm



7 (G-2)). If the difficulty of the test increases, the lower accuracy of GMRES (compared to the quadruple-precision routine used in LU QUAD Equil, LU QUAD, and the combined algorithm) leads to a greater number of Newton iterations, as shown in Fig. 7 for the MoMaS hard case. This point is confirmed for other cases (see Appendix 7 (G-4 to G-9)). For the Fe-Cr mineral case (see Appendix 7 (G-10)), only LU QUAD Equil and the combined algorithm can solve the problem. Other methods lead to non-convergence, due to overflow for the GMRES algorithm (overflow appears in the Newton algorithm and is not due to GMRES itself) and because LU QUAD and LU are unable to give an accurate descent direction.

Appendix 7 (G-11) shows the evolution of the relation between the norm of Y and the condition number of the Jacobian matrix during the minimization process. This figure is similar to Fig. 1, confirming the empirical relation (25). This relation cannot be used close to the solution, and the condition number tends to be a case-dependent limit for very low $\|Y\|$.

Nevertheless, the number of iterations is not the critical point. Because the time required by one iteration changes depending on the method used, we have to consider the total computation time. By plotting the total computation time required to solve each test case depending on the algorithm used (see Fig. 8), we can see that

- (i) LU QUAD Equil, as expected, is the slowest. Nevertheless, this method allows the convergence of the Newton-Raphson method for all test cases.
- (ii) LU QUAD is slightly faster. The difference between LU QUAD Equil and LU QUAD gives an indication of the time used for matrix equilibration. This time is greater for pyrite, MoMaS easy, pyrite mineral, MoMaS hard, and Fe-Cr mineral than for the other test cases.
- (iii) LU is fast when it leads to convergence, but this method results in a very weak Newton-Raphson algorithm.
- (iv) GMRES always results in the fastest Newton-Raphson algorithm. It has been shown (Fig. 7, Appendix 7 (G-8)) that the number of required iterations can be twice the number for other methods, but we show (Fig. 5) that the GMRES method is faster than the other methods.
- (v) The proposed combined algorithm leads to intermediate computing times, equivalent to those of LU QUAD Equil and LU QUAD, depending on the case.

According to our results, GMRES should be systematically used because it is fast and usually leads to convergence of the Newton-Raphson algorithm. The combined algorithm should be used for very high condition numbers or for recomputing a failed run.

5 Conclusion

In this work, we focus on the resolution of small linear systems generated using the Newton-Raphson algorithm to solve equilibrium chemistry problems. For the first time, we propose a study of the condition number of such linear systems and find that the range of values covered is unusually large. This characteristic leads to specific numerical problems, with matrices that are quite small (approximately 10×10) but very badly conditioned (up to 10^{100}). Ten different chemical systems are studied.

There is a strong linear relationship between the logarithm of the condition number of the matrix and the logarithm of the norm of the objective function. This factor can be exploited to create efficient algorithms. This relation is strictly an empirical one and is not valuable for low condition numbers.

A wide variety of linear solvers have been tested, and several direct and iterative solvers are selected. Some of these solvers are specific for a class of matrix, symmetric or positive definite, while others are generic. A preconditioning method (matrix equilibration) has also been tested to reduce the conditioning of the systems.

According to our selected test cases, only the LU and LU QUAD direct methods, the GMRES iterative method, and LU Equil and GMRES Equil preconditioned methods are sufficiently robust to solve all the tests.

According to the size of the chemical tests, the LU method is faster than the GMRES method. However, our results for the Fe-Cr mineral and Morel-Morgan cases show that GMRES is preferable for larger chemical systems (more than 10 components). Chemical systems with more than 10 components have not been frequently modelled in the past decade. However, the use of geochemical databases makes the construction of large geochemical systems easier, and the increase in computation capacities makes it possible. For very large geochemical systems, we recommend the GMRES method.

We also propose using the linear relationship between the condition number of the Jacobian matrix and the norm of the objective function to develop an efficient algorithm.

The classic LU method is not a good choice. Its weakness is its low robustness for challenging test cases. We recommend using the GMRES method, which is fast and usually leads to convergence of the Newton-Raphson algorithm. For very high condition numbers (more than 10^{100}), we recommend the most robust LU QUAD Equil method. When the Newton-Raphson method is sufficiently near the solution to decrease the condition number, the faster GMRES method can be used. By using the linear relationship between $\text{cond}(Z)$ and $\|Y\|$, the transition between the two methods can be achieved without computing the condition number (which is very expensive).

1043 This work explores a new research field by studying geo-
 1044 chemical computation from a condition number point of
 1045 view. We attempted to benchmark a wide variety of linear
 1046 solvers, but it was not possible to explore the flexibility
 1047 of all the tested solvers. This study will help us to eliminate
 1048 some solvers so that our future work can focus on the
 1049 most promising: LU, LU QUAD, GMRES, LU Equil, and
 1050 GMRES Equil. Some points for future exploration are as
 1051 follows:

- 1052 (i) We did not extensively test the robustness and the
 1053 efficiency of the Newton-Raphson algorithm. Further
 1054 work should examine the influence of the initial
 1055 Newton-Raphson guess to confirm our conclusions about
 1056 the high efficiency of the GMRES method.
 1057
- 1058 (ii) The accuracy of iterative methods depends on the
 1059 value of the convergence criterion (which we set to
 1060 10^{-8}) and on the method used to check the convergence
 1061 (we used the default method). Moreover, the efficiency
 1062 can vary depending on the initial guess provided by the
 1063 user. In this work, we used the easiest initial guess:
 1064 the residual for the tests from the Newton-Raphson
 1065 method and the previous Newton-Raphson step for the
 1066 test in a Newton-Raphson algorithm. We believe that
 1067 it is possible to make a better choice, markedly
 1068 enhancing the efficiency of the iterative methods.
 1069
- 1070 (iii) The GMRES method allows the use of left and/or
 1071 right preconditioners. These preconditioners can
 1072 increase the robustness, accuracy, and efficiency of
 1073 the method. More generally, several classes of
 1074 preconditioners that may reduce the condition number
 1075 of the linear system can be used [65, 66]. In this
 1076 work, we explored the use of one preconditioner:
 1077 matrix equilibration. However, other classes of
 1078 preconditioners may be more efficient.
- 1079 (iv) Previous works have addressed the use of methods
 1080 to solve geochemical equilibria other than the
 1081 Newton-Raphson method [17, 44, 49, 67]. It has
 1082 been shown [17] that an efficient algorithm can be
 1083 obtained by combining a zero-order method with
 1084 the Newton-Raphson approach.
- 1085 (v) The size of the chemical tests presented here is
 1086 representative of the sizes actually used in
 1087 environmental studies. We have shown that the
 1088 GMRES method may be efficient for large systems.
 1089 In anticipation of future needs, it may be useful
 1090 to test chemical systems larger than the Morel-
 1091 Morgan system.
- 1092 (vi) Part of the Newton minimization related to
 1093 very large condition numbers (far from the
 1094 solution) can be performed using *random* methods;
 1095 GMRES is efficient even though its descent direction
 1096 is not accurate

for high condition numbers. Some methods, such as
 simulated annealing and particle swarm optimization,
 could be used in future research.

These factors should be explored in light of the results
 presented in this study. We proposed a large set of
 chemical tests, a criterion to determine the difficulty
 of these tests (the condition number), and a panel
 of numerical methods that should be studied
 preferentially.

As a more general consideration, the reader should
 pay particular attention to the old Morel-Morgan
 test case and the more realistic pyrite test case.
 The Morel-Morgan test uses Fe^{2+} and Fe^{3+} , Co^{2+} and Co^{3+} ,
 and SO_4^{2-} and S^{2-} as components whereas the
 pyrite case uses O_2 , Fe^{2+} and SO_4^{2-} . The first
 studies on geochemical computation avoided redox
 problems. We show that redox problems lead to
 higher condition numbers because the stoichiometric
 coefficients and equilibrium constants cover a wider
 range. Several geochemical databases avoid the
 introduction of redox reactions. There is sometimes
 a good reason to not write redox reactions as
 equilibria (slow reaction rates, irreversible
 reactions) as done in Arora et al. [2]. However,
 the reason is sometimes numeric, and redox
 reactions are avoided because they lead to non-
 convergence.

We propose the use of quadruple-precision real
 for challenging chemical systems. In this work,
 the core of the geochemical code is conserved
 as double-precision real, and only the linear
 system tool is set as quadruple precision.
 Rewriting an entire geochemical code in a
 quadruple-precision format will result in robust
 code but at the cost of an important and
 rebarbative work as well as efficiency. In this
 stage of our research, we do not recommend
 such an effort because implementing LU
 decomposition using quadruple-precision real
 is very efficient, requiring only a minor
 modification of existing code and reducing
 the computation time.

Acknowledgments Hela Machat has been supported
 by a grant from the Tunisian Government. This
 work is supported by the BRGM-CNRS CUBICM
 project. We thank the reviewers for their helpful
 comments.

Compliance with Ethical Standards Compliance
 with ethical standards

Conflict of interests The authors declare that
 they have no conflict of interest.

Ethical approval This article does not contain
 any studies with human participants or animals
 performed by any of the authors.

Informed consent Informed consent was
 obtained from all individual participants included
 in the study.

Q10

Appendix 1

Q11142

1143 Morel's table of chemical test cases

1144

		X1	X2	X3	X4	X5	S	Log (K)
1	X1	1	0	0	0	0	0	0.00
2	X2	0	1	0	0	0	0	0.00
3	X3	0	0	1	0	0	0	0.00
4	X4	0	0	0	1	0	0	0.00
5	X5	0	0	0	0	1	0	0.00
6	C1	0	-1	0	0	0	0	-12.00
7	C2	0	1	1	0	0	0	0.00
8	C3	0	-1	0	1	0	0	0.00
9	C4	0	-4	1	3	0	0	-1.00
10	C5	0	4	3	1	0	0	35.00
11	C6	0	10	3	0	0	0	32.00
1145 12	C7	0	-8	0	2	0	0	-4.00
13	S	0	0	0	0	0	1	0.00
14	CS1	0	3	1	0	0	1	6.00
15	CS2	0	-3	0	1	0	2	-1.00
16	CP1 mineral	0	1	3	0	0	0	10.90
17	CP2 mineral	0	1	0	0	1	0	1.30
	Total (M)	0.3	0.3	0.3	2	0.3	10	
	X value	0.1	Variable	Variable	Variable	1.00E-03	1.00E-03	
	Min value		1.00E-1	1.00E-1	1.00E-1			
			5	5	5			
	Max value		1.00E-02	1.00E-02	1.00E-02			
	X initial value	0.1	1.00E-07	1.00E-07	1.00E-03	1.00E-03	1.00E-03	
	for Newton-Raphson iteration							

1146 **Appendix 2**

1147 Jacobian matrix properties

1148 B-1 Ratio of diagonal dominant matrices as a function of
1149 the condition number for the 10 chemical test cases

1150 B-2 Ratio of positive definite matrices as a function of
1151 the condition number for the 10 chemical test cases

C-21-C-27: Preconditioned methods, chemical cases without minerals 1162

C-28-C-30: Preconditioned methods, chemical cases with minerals 1164

Appendix 4 1166

Evolution for the relative residual on norm of the 16 selected methods as a function of the condition number for the 10 chemical test cases 1167

1152 **Appendix 3**

1153 Evolution of the failure ratio of the 16 selected methods as
1154 a function of the condition number for the 10 chemical test
1155 cases

1156 C-1-C-7: Direct methods, chemical cases without minerals

1157 C-8-C-10: Direct methods, chemical cases with minerals

1158 C-11-C-17: Iterative methods, chemical cases without
1159 minerals

1160 C-18-C-20: Iterative methods, chemical cases with
1161 minerals

D-1-D-7: Direct methods, chemical cases without minerals 1170

D-8-D-10: Direct methods, chemical cases with minerals 1171

D-11-D-17: Iterative methods, chemical cases without minerals 1172

D-18-D-20: Iterative methods, chemical cases with minerals 1174

D-31-D-37: Preconditioned methods, chemical cases without minerals 1175 Q12

D-38-D-40: Preconditioned methods, chemical cases with minerals 1177

Appendix 5

1179 Evolution for the angle between the reference and computed
 1180 solution of the 16 selected methods as a function of the
 1181 condition number for the 10 chemical test cases
 1182 E-1–E-7: Direct methods, chemical cases without minerals
 1183 E-8–E-10: Direct methods, chemical cases with minerals
 1184 E-11–E-17: Iterative methods, chemical cases without
 1185 minerals
 1186 E-18–E-20: Iterative methods, chemical cases with
 1187 minerals
 1188 E-21–E-27: Preconditioned methods, chemical cases
 1189 without minerals
 1190 E-28–E-30: Preconditioned methods, chemical cases
 1191 with minerals
 1192

Appendix 6

1193 Evolution for the computation time of the 16 selected meth-
 1194 ods as a function of the condition number for the 10
 1195 chemical test cases
 1196 F-1–F-7: Direct methods, chemical cases without minerals
 1197 F-8–F-10: Direct methods, chemical cases with minerals
 1198 F-11–F-17: Iterative methods, chemical cases without
 1199 minerals
 1200 F-18–F-20: Iterative methods, chemical cases with
 1201 minerals
 1202 F-21–F-27: Preconditioned methods, chemical cases
 1203 without minerals
 1204 F-28–F-30: Preconditioned methods, chemical cases
 1205 with minerals
 1206

Appendix 7

1207 Evolution of the Newton-Raphson residual as a function of
 1208 the number of iterations for the 10 chemical test cases and
 1209 the 5 tested algorithms
 1210 G-1–G-7: Chemical test without mineral
 1211 G-8–G-10: Chemical test with minerals
 1212 G-11: Evolution of the relation between $\|Y\|$ and the
 1213 condition number of the Jacobian matrix Z during
 1214 minimization
 1215

References

1216 1. Walter, A.L. et al.: Modeling of multicomponent reactive transport
 1217 in groundwater. 2. Metal mobility in aquifers impacted by acidic
 1218 mine tailings discharge. *Water Resour. Res.* **30**(11), 3149–3158
 1219 (1994)
 1220

2. Arora, B. et al.: A reactive transport benchmark on heavy metal 1221
 cycling in lake sediments *Computational Geosciences* (2014) 1222
 3. De Windt, L., Leclercq, S., Van der Lee, J.: Assessing the durability 1223
 of nuclear glass with respect to silica controlling processes in 1224
 a clayey underground disposal. In: 29th International Symposium 1225
 on the Scientific Basis for Nuclear Waste Management XXIX. 1226
 Materials Research Society Symposium Proceedings, Ghent; Bel- 1227
 gium (2005) 1228
 4. Hoteit, H., Ackerer, P., Mose, R.: Nuclear waste disposal simula- 1229
 tions: Couplex test cases. *Comput. Geosci.* **8**(2), 99–124 (2004) 1230
 5. Tompson, A.F.B., et al.: On the evaluation of groundwater contam- 1231
 ination from underground nuclear tests. *Environ. Geol.* **42**(2-3), 1232
 235–247 (2002) 1233
 6. Andre, L., et al.: Numerical modeling of fluid-rock chemical 1234
 interactions at the supercritical CO₂-liquid interface during CO₂ 1235
 injection into a carbonate reservoir, the Dogger aquifer (Paris 1236
 Basin, France). *Energy Convers. Manag.* **48**(6), 1782–1797 (2007) 1237
 7. Kang, Q., et al.: Pore scale modeling of reactive transport involved 1238
 in geologic CO₂ sequestration. *Transp. Porous Media* **82**(1), 197– 1239
 213 (2010) 1240
 8. Navarre-Sitchler, A.K., et al.: Elucidating geochemical response 1241
 of shallow heterogeneous aquifers to CO₂ leakage using high- 1242
 performance computing: implications for monitoring of CO₂ 1243
 sequestration. *Adv. Water Resour.* **53**(0), 45–55 (2013) 1244
 9. Pruess, K. et al.: Code intercomparison builds confidence in 1245
 numerical simulation models for geologic disposal of CO₂. 1246
Energy **29**(9-10), 1431–1444 (2004) 1247
 10. Regnault, O., et al.: Etude experimentale de la reactivite du CO₂ 1248
 supercritique vis-a-vis de phases minerales pures. Implications 1249
 pour la sequestration geologique de CO₂. *Compt. Rendus Geosci.* 1250
337(15), 1331–1339 (2005) 1251
 11. Valocchi, A.J., Street, R.L., Roberts, P.V.: Transport of ion- 1252
 exchanging solutes in groundwater: chromatographic theory and 1253
 field simulation. *Water Resour. Res.* **17**, 1517–1527 (1981) 1254
 12. Lichtner, P.C.: Continuum model for simultaneous chemical reac- 1255
 tions and mass transport in hydrothermal systems. *Geochim.* 1256
Cosmochim. Acta **49**(3), 779–800 (1985) 1257
 13. Appelo, C.A.J.: Hydrogeochemical transport modelling. *Proceed.* 1258
Inf.—Comm. Hydrol. Res. TNO **43**, 81–104 (1990) 1259
 14. Yeh, G.T., Tripathi, V.S.: A critical evaluation of recent develop- 1260
 ments in hydrogeochemical transport models of reactive multi- 1261
 chemical components. *Water Resour. Res.* **25**, 93–108 (1989) 1262
 15. Carrayrou, J. et al.: Comparison of numerical methods for sim- 1263
 ulating strongly nonlinear and heterogeneous reactive transport 1264
 problems—the MoMaS benchmark case. *Computational Geo-* 1265
sciences **14**(3), 483–502 (2010) 1266
 16. Hammond, G.E., Valocchi, A.J., Lichtner, P.C.: Modeling mul- 1267
 ticomponent reactive transport on parallel computers using 1268
 Jacobian-Free Newton Krylov with operator-split preconditioning. 1269
 In: Hassanzadeh, S.M. (ed.) *Developments in water science, com-* 1270
putational methods in water resources, Proceedings of the XIVth 1271
International Conference on Computational Methods in Water 1272
Resources (CMWR XIV), pp. 727–734. Elsevier (2002) 1273
 17. Carrayrou, J., Mosé, R., Behra, P.: New efficient algorithm for 1274
 solving thermodynamic chemistry. *AIChE J.* **48**(4), 894–904 1275
 (2002) 1276
 18. Amir, L., Kern, M.: A global method for coupling transport with 1277
 chemistry in heterogeneous porous media. *Comput. Geosci.* **14**(3), 1278
 465–481 (2010) 1279
 19. Quarteroni, A., Sacco, R., Saleri, F.: Numerical mathematics. In: 1280
 Marsden, J.E., Sirovich, L., Antman, S.S. (eds.) *Texts in Applied* 1281
Mathematics. 2nd edn. Springer, Heidelberg (2007) 1282
 20. Axelsson, O., et al.: Direct solution and incomplete factoriza- 1283
 tion preconditioned conjugate gradient methods. Comparison of 1284

1285 algebraic solution methods on a set of benchmark problems in linear
 1286 elasticity, in STW report. 2000, Department of Mathematics,
 1287 Catholic University of Nijmegen: Nijmegen, The Netherlands. pp.
 1288 1–36

1289 21. Barrett, R., Berry, M., Chan, T.F., Demmel, J., Donato, J., Don-
 1290 garra, J., Eijkhout, V., Pozo, R., Romine, C., Van Der Vorst, H.
 1291 Templates for the solution of linear systems: building blocks for
 1292 iterative methods, 2nd edn. SIAM, Philadelphia (1994)

1293 22. Gould, N.I.M., Hu, Y., Scott, J.A.: A numerical evaluation of
 1294 sparse direct solvers for the solution of large sparse, symmetric lin-
 1295 ear systems of equations. 2005, Council for the Central Laboratory
 1296 of the Research Councils

1297 23. Allaire, G., Kaber, S.M. In: Marsden, J.E., Sirovich, L., Antman,
 1298 S.S. (eds.): Numerical linear algebra. Texts in applied mathemat-
 1299 ics. Springer, New York (2008)

1300 24. Baldwin, C., et al.: Iterative linear solvers in a 2D radiation-
 1301 hydrodynamics code: methods and performance. *J. Comput. Phys.*
 1302 **154**(1), 1–40 (1999)

1303 25. Chao B.T., L.H.L., Scott, E.J.: On the solution of ill-conditioned,
 1304 simultaneous, linear, algebraic equations by machine computation,
 1305 in *University of Illinois Bulletin*. 1961, University of Illinois

1306 26. Hadjidimos, A.: Successive overrelaxation (SOR) and related
 1307 methods. *J. Comput. Appl. Math.* **123**(1-2), 177–199 (2000)

1308 27. Kalambi, I.B.: A comparison of three iterative methods for the
 1309 solution of linear equations. *J. Appl. Sci. Environ. Manag.* **12**(4),
 1310 53–55 (2008)

1311 28. Klisinski, M., Runesson, K.: Improved symmetric and non-
 1312 symmetric solvers for FE calculations. *Adv. Eng. Softw.* **18**(1),
 1313 41–51 (1993)

1314 29. Schenk, O., Gartner, K.: Solving unsymmetric sparse systems of
 1315 linear equations with PARDISO. *Fut. Gener. Comput. Syst.* **20**(3),
 1316 475–487 (2004)

1317 30. Xue, X.J. et al.: A direct algorithm for solving ill-conditioned
 1318 linear algebraic systems. *JCPDS-Int. Centre Diffract. Data Adv.*
 1319 *X-ray Anal.* **42**, 629–633 (2000)

1320 31. Pyzara, A., Bylina, B., Bylina, J.: The influence of a matrix
 1321 condition number on iterative methods' convergence (2011)

1322 32. Hoffmann, J., Kräsutle, S., Knabner, P.: A parallel global-implicit
 1323 2-D solver for reactive transport problems in porous media based
 1324 on a reduction scheme and its application to the MoMaS bench-
 1325 mark problem. *Comput. Geosci.* **14**(3), 421–433 (2010)

1326 33. Soleymani, F.: A new method for solving ill-conditioned linear
 1327 systems. *Opuscula Math.* **33**(2), 337–344 (2013)

1328 34. Morel, F., Morgan, J.: A numerical method for computing equi-
 1329 libria in aqueous chemical systems. *Environ. Sci. Technol.* **6**(1),
 1330 58–67 (1972)

1331 35. Morel, F.M.M.: Principles of aquatic chemistry. Wiley Inter-
 1332 science, New York (1983)

1333 36. De Windt, L. et al.: Intercomparison of reactive transport models
 1334 applied to UO₂ oxidative dissolution and uranium migration. *J.*
 1335 *Contam. Hydrol.* **61**(1-4), 303–312 (2003)

1336 37. Jauzein, M. et al.: A flexible computer code for modelling
 1337 transport in porous media: impact. *Geoderma* **44**(2–3), 95–
 1338 113 (1989)

1339 38. Parkhurst, D.L., Appelo, C.A.J.: User's guide to PHREEQC
 1340 (version 2)—a computer program for speciation, batch-reaction,
 1341 one-dimensional transport, and inverse geochemical calculations.
 1342 *Water Resour. Invest.*, Editor. 1999: Denver. p. 312

1343 39. Van der Lee, J.: CHESS another speciation and surface complex-
 1344 ation computer code. E.d.M.d. Paris, Editor. 1993: Fontainebleau.
 1345 p. 85

1346 40. Westall, J.C.: MICROQL: a chemical equilibrium program in
 1347 BASIC. Computation of adsorption equilibria in BASIC. S.F.I.o.T.
 1348 EAWAG, Editor. 1979: Dübendorf. p. 42

41. Westall, J.C.: FITEQL ver. 2.1. 1982: Corvallis 1349

42. Westall, J.C., Zachary, J.L., Morel, F.M.M.: MINEQL: a computer 1350
 program for the calculation of chemical equilibrium composition 1351
 of aqueous system. R.M.P. Laboratory, Editor. 1976: Cambridge. 1352
 p. 91 1353

43. Walter, L.J., Wolery, T.J.: A monotone-sequences algorithm and 1354
 FORTRAN IV program for calculation of equilibrium distribu- 1355
 tions of chemical species. *Comput. Geosci.* **1**, 57–63 (1975) 1356

44. Wigley, T.M.L.: WATSPEC: a computer program for determining 1357
 the equilibrium speciation of aqueous solutions. B.G.R.G. Tech. 1358
 Bull., Editor. 1977. p. 49 1359

45. Jennings, A.A., Kirkner, D.J., Theis, T.L.: Multicomponent equi- 1360
 librium chemistry in groundwater quality models. *Water Resour.* 1361
Res. **18**, 1089–1096 (1982) 1362

46. Cederberg, A., Street, R.L., Leckie, J.O.: A groundwater mass 1363
 transport and equilibrium chemistry model for multicomponent 1364
 systems. *Water Resour. Res.* **21**, 1095–1104 (1985) 1365

47. Yeh, G.T., Tripathi, V.S.: A model for simulating transport of reac- 1366
 tive multispecies components: model development and demonstra- 1367
 tion. *Water Resour. Res.* **27**(12), 3075–3094 (1991) 1368

48. Carrayrou, J.: Looking for some reference solutions for the 1369
 reactive transport benchmark of MoMaS with SPECY. *Comput.* 1370
Geosci. **14**(3), 393–403 (2010) 1371

49. Brassard, P., Bodurtha, P.: A feasible set for chemical speciation 1372
 problems. *Comput. Geosci.* **26**(3), 277–291 (2000) 1373

50. Carrayrou, J., Kern, M., Knabner, P.: Reactive transport bench- 1374
 mark of MoMaS. *Comput. Geosci.* **14**(3), 385–392 (2010) 1375

51. Fendorf, S.E., Li, G.: Kinetics of chromate reduction by ferrous 1376
 iron. *Environ. Sci. Technol.* **30**(5), 1614–1617 (1996) 1377

52. Chilakapati, A. et al.: Groundwater flow, multicomponent trans- 1378
 port and biogeochemistry: development and application of a 1379
 coupled process model. *J. Contam. Hydrol.* **43**(3-4), 303–325 1380
 (2000) 1381

53. Knight, P., Ruiz, D., Ucar, B.: A symmetry preserving algorithm 1382
 for matrix scaling. *SIAM J. Matrix Anal. Appl.* **35**(3), 931–955 1383
 (2014) 1384

54. Golub, H.V., Van Loan, C.F.: Matrix computations. 3rd ed. The 1385
 Johns Hopkins University Press, Baltimore (1996) 1386

55. Davis, T.A., Duff, I.S.: A combined unifrontal/multifrontal 1387
 method for unsymmetric sparse matrices. *ACM Trans. Math.* 1388
Softw. **25**(1), 1–20 (1999) 1389

56. Woźnicki, Z.: On performance of SOR method for solving non- 1390
 symmetric linear systems. *J. Comput. Appl. Math.* **137**(1), 145– 1391
 176 (2001) 1392

57. Saad, Y., Van Der Vorst, H.A.: Iterative solution of linear systems 1393
 in the 20th century. *J. Comput. Appl. Math.* **123**(1-2), 1–33 (2000) 1394

58. Diersch, H.J.G.: FEFLOW reference manual. DHI-WASY GmbH, 1395
 Berlin (2009) 1396

59. Van der Lee, J., et al.: Presentation and application of the reactive 1397
 transport code HYTEC. In: Hassanizadeh, S.M. (ed.) Develop- 1398
 ments in Water Science, Computational Methods in Water 1399
 Resources, Proceedings of the XIVth International Conference 1400
 on Computational Methods in Water Resources (CMWR XIV), 1401
 pp. 599–606. Elsevier (2002) 1402

60. Press, W.H., S.A.T., Vetterling, W.T., Flannery, B.P. Numerical 1403
 recipes in FORTRAN: the art of scientific computation, 2nd edn., 1404
 pp. 123–124. Cambridge University Press, New Yor (1992) 1405

61. The Linear Algebra Package (LAPACK) can be obtained free of 1406
 charge from the address listed here: <http://www.netlib.org/lapack> 1407

62. Kincaid, D., Cheney, W. Numerical analysis: mathematics of 1408
 scientific computing, 3rd edn. American Mathematical Society 1409
 (2002) 1410

63. HSL: A collection of Fortran codes for large scale scientific 1411
 computation. <http://www.hsl.rl.ac.uk> (2013) 1412

- 1413 64. Chapter 8 Systems of nonlinear equations. In: Studies in com- 1418
1414 putational mathematics, Claude, B. Editor. 1997, Elsevier. pp. 1419
1415 287–336
1416 65. Soleymani, F.: A rapid numerical algorithm to compute matrix 1420
1417 inversion. *Int. J. Math. Math. Sci.* **2012** (2012) 1421
1422
66. Soleymani, F.: On a fast iterative method for approximate inverse 1418
of matrices. *Commun. Korean Math. Soc.* **28**(2), 407–418 (2013) 1419
67. Morin, K.A.: Simplified explanations and examples of comput- 1420
erized methods for calculating chemical equilibrium in water. 1421
Comput. Geosci. **11**, 409–416 (1985) 1422

UNCORRECTED
PROOF

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES:

- Q1. Please check the provided city (Jendouba) and country (Tunisia) names in affiliation 2 if correct.
- Q2. Please check if the abbreviation "LAPACK" is defined correctly. Otherwise, please provide the correct expansion.
- Q3. Please check all equation citations if correctly captured or presented.
- Q4. The section citation "Properties of the Jacobian matrix" found in the sentence starting "This formulation has some weaknesses" was changed to "Section 3.1" as per numbering style. Please check if correct.
- Q5. As per instruction, if there are more than one appendix, they should be designated with numbers 1, 2, 3, etc. With this regard, the appendices were renumbered accordingly and their citations were modified. Please advise if the labels designated for figures found in appendices B to G (as originally labelled) should also be modified accordingly to ensure consistency of presentation.
- Q6. Missing citation for Table 1 was inserted in this sentence. Please check if appropriate. Otherwise, please provide the location of where to insert the citation/s in the main body of the text. Note that the order of main citations of tables in the text must be sequential.
- Q7. Please check the captured caption of Figure 4 if correct.
- Q8. Missing citation for Figure 4 was inserted in this sentence. Please check if appropriate. Otherwise, please provide the location of where to insert the citation/s in the main body of the text. Note that the order of main citations of figures in the text must be sequential.
- Q9. The appendix citations found in the last two paragraphs before the "Proposal of a new algorithm" section were all changed to "Appendix 6" so that appendix citations are in sequential order. Also, the discussion in last two paragraphs fits to the description of Appendix 6. Please check if correct and amend if necessary.
- Q10. Please confirm if the "Informed consent" statement is indeed applicable to this article and should be retained.
- Q11. Please check if the content of Appendices 1 to 7 is correctly captured or presented.
- Q12. Please check if the labels "D-31 to D-40" are indeed correct in the artwork of Appendix 4, and not "D-21 to D-30", respectively.