

This is the preprint of the contribution published as:

Buchwald, J., Chaudhry, A.A., Yoshioka, K., Kolditz, O., Attinger, S., Nagel, T. (2020):
DoE-based history matching for probabilistic uncertainty quantification of thermo-hydro-
mechanical processes around heat sources in clay rocks
Int. J. Rock Mech. Min. Sci. **134** , art. 104481

The publisher's version is available at:

<http://dx.doi.org/10.1016/j.ijrmms.2020.104481>

1 DoE-based history matching for probabilistic uncertainty quantification of
2 thermo-hydro-mechanical processes around heat sources in clay rocks

3 J. Buchwald^{a,b,*}, A. A. Chaudhry^b, K. Yoshioka^a, O. Kolditz^{a,c}, S. Attinger^a, T. Nagel^b

4 ^a*Helmholtz Centre for Environmental Research – UFZ, Permoserstraße 15, 04318 Leipzig, Germany*

5 ^b*Chair of Soil Mechanics and Foundation Engineering, Geotechnical Institute, Technische Universität Bergakademie
6 Freiberg, Freiberg, Germany*

7 ^c*Applied Environmental Systems Analysis, Technische Universität Dresden, Germany*

8 **Abstract**

In the context of geotechnical and geological barriers, a thorough analysis of uncertainty and sensitivity is a crucial aspect of any physics-based performance assessment. While experimental data are scarce in actual waste repositories, large-scale experiments in underground research laboratories (URLs) provide such data that can be used to not only qualify THMC process models but also uncertainty assessment methodologies. In this paper, we adopt a Design of Experiments (DoE)-based history matching workflow – an approach popular in the oil and gas industry – and scrutinize its applicability for multiphysical analyses of nuclear waste disposal-related processes using synthetic experimental data. Based on an analytical solution of a coupled thermo-hydro-mechanical (THM) problem of a heat source embedded in a fluid-saturated porous medium mimicking a disposal cell in an argillaceous host formation, we discuss the adaptability of the workflow as a way to address parameter and model uncertainties for barrier integrity assessment. We thereby put particular focus on the relative importance of providing defined input parameter distributions for quantities generally afflicted with epistemic uncertainty and the constraints imposed by experimental (URL) or monitoring (repository) data. We found that once constraining data is available, the particular a priori distribution plays only a minor role for the outcome, such that we can conclude that the often unknown distributions can be substituted by uniform priors under such conditions. However, detailed knowledge of parameter distributions can increase the efficiency of the workflow significantly. We conclude that the presented workflow is particularly suitable for performing uncertainty quantification and sensitivity analysis for geotechnical applications where monitoring or other experimental data are available, as it allows us to deal with models of great complexity, epistemic uncertainty and it incorporates canonically to use of measured data in order

to reduce uncertainty.

9 *Keywords:* Design of Experiments, history matching, thermo-hydro-mechanical, radioactive
10 waste, geological repository, uncertainty analysis, uncertainty quantification, sensitivity analysis,
11 OpenGeoSys
12 *2000 MSC:* 62K20
13 *2000 MSC:* 49Q12
14 *2000 MSC:* 60-08
15 *2000 MSC:* 60G15

16 **1. Introduction**

17 One of the objectives in the design of nuclear waste repositories is to provide solutions that are
18 robust and simple in the sense that their future evolution can be subjected to predictive analyses
19 based on the current physical understanding of the involved processes. Nevertheless, in conjunction
20 with the large spatial and temporal scales relevant to system evolution, any predictive analysis of
21 thermo-hydro-mechanical processes around radioactive waste repositories retains a considerable de-
22 gree of uncertainty. Safety regulations, therefore, require implementors to address this uncertainty
23 at different levels.

24 The focus of this contribution is uncertainty quantification for the use in performance assess-
25 ment of radioactive waste repositories and their components, which remains a thoroughly challeng-
26 ing task. Commonly, one distinguishes between uncertainty analysis/quantification and sensitivity
27 analysis. The former is related to the determination of the overall uncertainty of a model system
28 in terms of its input, whereas the latter refers to evaluating the relative contribution of each input
29 (cf. e.g., [1]) to this uncertainty. Both aspects are naturally linked and addressed in this work.
30 Dependent on the host rock in which nuclear waste is to be emplaced, different thermo-hydro-
31 mechanical-chemical (THMC) processes need to be analyzed that might influence the transport of
32 radionuclides in the geological disposal system [2, 3, 4]. In a repository for high-level waste, some
33 of the most significant effects during the post-closure phase are triggered by the waste packages'
34 decay heat, causing major changes in the physical properties of the host-rock in the near-field and
35 driving the system away from its former equilibrium state.

*Corresponding author

36 Many other disturbances to the natural state are caused by the excavation itself and dominate
37 the early stages of a repository evolution. The excavation-induced changes of the HM boundary
38 conditions and chemical equilibria in the near-field also represent initial conditions for the post-
39 closure phase. Such effects are neglected here. A detailed overview and description of processes
40 that play an important role can be found elsewhere [5, 3, 4].

41 J. C. Helton [6] provided a general overview of different uncertainty and sensitivity analysis
42 techniques applicable to a broad set of problems. Most of the outlined techniques remain perti-
43 nent, while some new methods emerged (e.g., SSFEM [7, 8, 9], random set theory [10, 8], or new
44 approaches to address global sensitivity [11, 12]) or existing methods have been developed further
45 since the early nineties (response-surface methods [13, 14] and Monte-Carlo methods [15]). The
46 choice of a particular method depends—aside from conceptual aspects—to a large extent on the
47 specific problem, its complexity, and the available computing power.

48 Indeed, when performing distribution sampling, the resulting response distributions will heavily
49 depend on the a priori distributions of the input variables. In the context of studying radioactive
50 waste disposal, this might be crucial. Commonly, we distinguish between *aleatory* and *epistemic*
51 uncertainty. Whereas aleatory uncertainty is due to a random process or a stochastic variability
52 of a phenomenon, epistemic uncertainty represents insufficient knowledge about a parameter (cf.
53 e.g., [16]). Here, we are dealing mostly with epistemic uncertainty because most parameters vary
54 intricately in space and time, meaning that it would be, in most cases, an onerous task to precisely
55 measure them. It is a matter of on-going debate, whether it is possible in all cases to adequately
56 describe these kinds of uncertainties in terms of classical probability theory [17, 18, 10, 19].

57 In our work, by model error, we are referring to structural discrepancies coming from a lack of
58 knowledge of the underlying physics for the problem at hand. This also emphasizes the need to
59 validate a computer model (this includes not only the governing equations and numerics but also
60 the geometry, scales that need to be considered, and also its boundaries). In order to judge the
61 practical significance of a validation, the model’s uncertainty needs to be quantified. Therefore, an
62 approach to uncertainty quantification that incorporates experimental data that can be directly
63 compared to model predictions would be very appealing. Aside from parameter uncertainties and
64 model inadequacy, other sources of uncertainties constitute numerical inaccuracies and observation
65 errors that are not tackled by our investigation as they can be subsumed by the others or are of

66 low relevance in most cases.

67 The software employed for performance assessment and coupled process simulations of geotech-
68 nical and geological barriers in repositories in various host rocks has undergone significant develop-
69 ment over the past decades both in terms of physical representation and computational efficiency
70 [20, 21, 22, 23, 24]. This development, however, has been largely disconnected from the develop-
71 ments in the UQ [25, 26] community, and very few links exist [27]. Thus, while improvements in
72 hardware and software over the past decade have provided us with modeling tools that enable the
73 modeling of radioactive waste repository sites in increasing detail, a direct Monte-Carlo approach
74 would still pose enormous demands in terms of computing power and time. Response surface,
75 proxy, or surrogate model approaches are one avenue for keeping the analyses tractable.

76 While in past discussions of parameter and model uncertainties in radioactive waste repositories,
77 the focus has been mainly on transport phenomena [28, 16, 6, 29], a future challenge remains to put
78 the entire coupled system (i.e. THMC processes) with all its relevant uncertainties under scrutiny
79 as only a few stochastic case studies exist to the present day [30, 31, 32].

80 Data to calibrate and validate models against require a monitoring program to be installed
81 as part of repository construction. For nuclear waste disposal research, it discloses another im-
82 portant application area: model development and validation using large-scale in-situ experiments
83 in underground research laboratories cannot merely be about chasing the 'best match'. Instead, a
84 meaningful model development and validation endeavour needs to take both experimental and mod-
85 eling uncertainty into account. Global sensitivity analysis (GSA) and uncertainty quantification
86 (UQ) can help not only in obtaining a better physical understanding but also in the identification
87 of meaningful target validation corridors for modelers to aim at.

88 This article focuses on the applicability of the DoE-based history matching workflow to the
89 multiphysical modeling of radioactive waste repository sites and components. In lieu of validating
90 the approach based on ongoing in-situ experiments and envisioning the implementation of a certain
91 monitoring phase of a future repository, our objective is to calibrate the THM model against
92 continuously monitored in-situ data and subsequently to perform a probabilistic predictive analysis,
93 i.e., a forecast. To pave the way for such analyses, this paper evaluates the workflow based on
94 synthetic experimental data.

95 **2. The underlying multiphysical (THM) problem**

The DoE-based history matching workflow is applied to a coupled thermo-hydro-mechanical (THM) model of a heat source embedded in an isotropic fluid-saturated porous medium. Although clay-rock typically exhibits thermal, hydraulic, and mechanical anisotropy, this was ignored for the purposes of this study, which used synthetic experimental data derived from an isotropic model. The same workflow, as described below, can be applied to more general settings without the need of further changes. The model can be formulated in terms of the three primary variables temperature, pore pressure, and displacement in the balance equations of mass, momentum, and energy complemented by the constitutive relationships of the fluid and solid phases as well as their interaction in the context of porous-media mechanics [33, 34]. As the heat source causes an increase in local temperatures, solids and fluids expand, creating pore pressure and effective stress variations. The emerging pressure gradient causes the fluid to flow away from the heat source, resulting thereby in a dissipation of the pore pressure in a thermally induced consolidation process. The corresponding equations of the linear problem can be written in terms of a thermal, hydraulic and mechanical part that are coupled to each other. The thermal part is described in terms of the energy balance equation which reads (for a brief nomenclature, see Tab. 1)

$$m\dot{T} + \rho_w c_w T_{,i} v_i - (KT_{,i})_{,i} = q_T \quad (1)$$

where q_T is a heat source per unit volume and, m (volumetric heat capacity), K (heat conductivity) and v_i (Darcy velocity) are given as

$$m = \phi \rho_w c_w + (1 - \phi) \rho_s c_s \quad (2)$$

$$K = \phi K_w + (1 - \phi) K_s \quad (3)$$

$$v_i = -\frac{k_s}{\mu} (p_{,i} - \rho_w g_i) \quad (4)$$

The the mass balance equation describes the hydraulic part including couplings and is given by

$$\beta \dot{p} - a_u \dot{T} + \alpha_B \dot{u}_{i,i} + v_{i,i} = q_H \quad (5)$$

where q_H is the source term for the fluid while a_u is given as

$$a_u = \phi a_w + (1 - \phi) a_s \quad (6)$$

The mechanical part can be derived from the momentum balance equations and reads

$$\sigma_{ij,j} + \rho g_i = 0 \quad (7)$$

where $\rho = \phi\rho_w + (1 - \phi)\rho_s$ and σ_{ij} is the total stress which is given by

$$\sigma_{ij} = \sigma'_{ij} - \alpha_{BP}p\delta_{ij} \quad (8)$$

where δ_{ij} refers to Kronecker delta and σ'_{ij} is the effective stress tensor which is given as

$$\sigma'_{ij} = C_{ijkl} \left(\epsilon_{kl} - \frac{1}{3}a'\Delta T\delta_{kl} \right) \quad (9)$$

Where C_{ijkl} and ϵ_{kl} are the stiffness and strain tensors, respectively. For isotropic case, the above equation can be rewritten as

$$\sigma'_{ij} = 2G\epsilon_{ij} + \lambda\epsilon_{kk}\delta_{ij} - b'\Delta T\delta_{ij} \quad (10)$$

where

$$b' = \left(\lambda + \frac{2G}{3} \right) a'. \quad (11)$$

96 This specific problem can be solved analytically under few simplifying assumptions. The so-
 97 lution and its comparison with a corresponding numerical model can be found elsewhere [35, 36].
 98 The model can be understood as a simplified version of a single disposal cell filled with radioactive
 99 material and emitting decay heat emplaced in an underground repository in a low-permeability
 100 host rock such as clay rock. In our workflow, the analytical model (further denoted as *AM*) is used
 101 on the one hand for the generation of synthetic experimental data, and on the other for the major
 102 part of the system analysis, i.e., for proxy generation (the kriging proxy model is further referred
 103 to as *PM*) as well as for history match selection/parameter identification.

104 In order to demonstrate the role of the forecast, i.e. a predictive analysis, in the workflow, we
 105 increased the power of the heat-flux in a step-wise manner, as done in actual heater experiments
 106 [37]. As a consequence, the resulting forecast response curves are non-trivial and predict behavior
 107 that goes beyond the calibration domain of the model. As the analytical solution is not capable
 108 of reflecting this change in heat power, the forecast uses the finite element method instead (the
 109 numerical model is further denoted as *NM* in the text). For this purpose, a two-dimensional model
 110 domain with axial symmetry representing a half-space of the spherically symmetric problem was

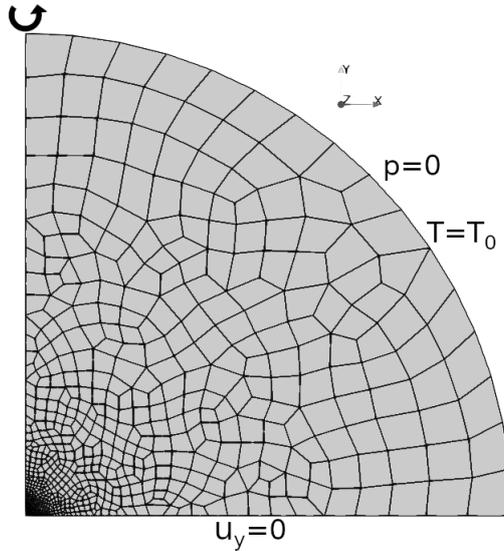


Figure 1: Mesh used for the numerical model with the applied boundary conditions.

111 created. Alongside the inner symmetry boundaries, the axis-normal displacements are set to zero,
 112 while at the outer boundaries the pore pressure is required to vanish¹, and the temperature is fixed
 113 to its initial value (293.15 K) along the outer boundaries. The mesh together with a summary of
 114 the applied boundary conditions is depicted in Fig. 1.

115 Further computational details are given in [36]. In Tab. 1, we present a list of parameters along
 116 with their lower and upper limits that will be used in the uncertainty and sensitivity analyses. The
 117 corresponding probability density functions are displayed in Fig. 1 of the electronic supplementary
 118 information (SI).

119 **3. History matching (HM) and uncertainty quantification (UQ)**

120 *3.1. General workflow*

121 The applied approach, known as experimental design (DoE)-based history matching, is closely
 122 related to the above-mentioned response-surface methods and Monte-Carlo sampling and has been
 123 applied to production forecasts in the oil and gas industry [13, 38]. To the best of our knowledge it
 124 has not been applied to problems in radioactive waste management. This is likely because usually,
 125 there is no historical data to match and most of the analyses are purely predictive. However,
 126 underground research laboratories (URLs) provide experimental data sets that allow us to explore

¹We depart from vanishing initial stresses and pore pressures due to the linearity in the constitutive models used. Stresses and pore pressures thus represent increments rather than absolute values.

Table 1: Used bounds of material parameters for the analytical (*AM*) and numerical model (*NM*). More details given in Tab. 1, SI.

Parameter	symbol	low	best	high	unit
Young’s modulus	E	$2.1 \cdot 10^9$	$2.7 \cdot 10^9$	$3.5 \cdot 10^9$	Pa
Poisson’s ratio	ν	0.28	0.33	0.38	-
Vol. thermal expansion coefficient of the solid	$a_s = a'$	$1.5 \cdot 10^{-6}$	$4.2 \cdot 10^{-6}$	$1 \cdot 10^{-5}$	K^{-1}
Vol. thermal expansion coefficient of water	a_w	$1.695 \cdot 10^{-4}$	$3.98 \cdot 10^{-4}$	$5.63 \cdot 10^{-4}$	K^{-1}
Porosity	ϕ	0.14	0.183	0.247	-
Water density	ρ_w	979.4736	991.46	998.767	$kg\ m^{-3}$
Solid grain density	ρ_s	2700.0	2768.5	2872.0	$kg\ m^{-3}$
Specific isobaric heat capacity of water	c_w	3941.38	4065.12	4167.71	$J\ kg^{-1}K^{-1}$
Specific isobaric heat capacity of the solid	c_s	760.0	860.0	960.0	$J\ kg^{-1}K^{-1}$
Heat conductivity of water	K_w	0.592015	0.63122	0.657	$W\ m^{-1}K^{-1}$
Heat conductivity of the solid	K_s	1.0	1.7	3.1	$W\ m^{-1}K^{-1}$
Dynamic viscosity of water	μ	$4.237 \cdot 10^{-4}$	0.000633	0.0011	Pa s
Intrinsic permeability	k_s	$2 \cdot 10^{-20}$	$3 \cdot 10^{-20}$	$2 \cdot 10^{-19}$	m^2
Initial temperature	T_0	292.15	293.15	294.15	K

127 this avenue. Therefore, an “ansatz” containing the parallel analysis of modeling and experimental
128 data is the ideal choice, as it allows us to link model calibration with the analysis of parameter
129 uncertainties. The applied approach is very closely linked to Bayesian history matching [39].
130 However, in contrast to the latter, the posterior function is taken from filtering the response of
131 the direct sampled priors, whereas Bayesian approaches typically use a likelihood function (derived
132 from the history-match error) to obtain Markov Chain Monte Carlo estimates.

133 Aside from the URL perspective, repository concepts including monitoring activities at least
134 for the early post-closure phase constitute a potential basis for analyses incorporating history
135 matching. Thus, if deemed suitable, this workflow can be of potential relevance for the design and
136 monitoring of future repository systems.

137 Note in passing, that similar considerations apply to other geotechnologies.

138 Following the steps as presented by [38], we want to highlight some features that appear par-
139 ticularly interesting for the class of problems at hand as summed up in the previous section:

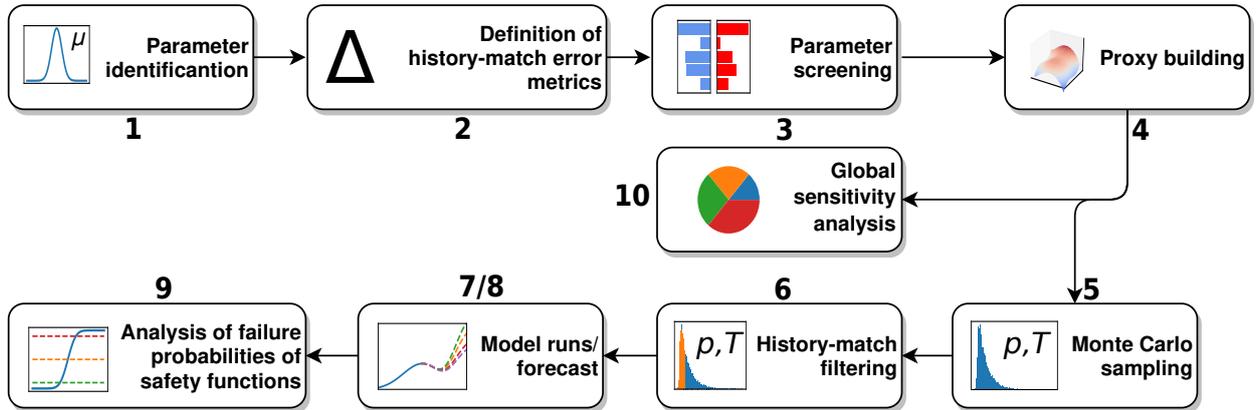


Figure 2: Schematic sketch of the workflow.

- 140 – history matching canonically incorporates experimental data enabling us to calibrate the
- 141 numerical model
- 142 – the use of proxy models for the history match error allows for direct Monte-Carlo sampling
- 143 with a statistically sufficient number of samples while keeping the computational burden
- 144 manageable
- 145 – initial parameter screening prior to proxy building makes it possible to neglect insignificant
- 146 parameters
- 147 – in contrast to many other more specialized uncertainty quantification methods in finite ele-
- 148 ment modeling, it is generally applicable to non-linear and coupled problems

149 In Fig. 2, we present a sketch in which we transferred the workflow to our purposes of uncertainty
 150 quantification and sensitivity analysis.

- 151 1. The first step is devoted to problem framing. In order to assess the applicability of this
- 152 approach to coupled THM models while keeping the computational burden minimal, we
- 153 employ an analytical solution that retains most of the essential primary THM couplings for
- 154 the model evaluations as well as for synthetic data generation. Another issue we want to
- 155 address in this paper is the sensitivity of this methodology to particular input parameter
- 156 distributions given only lower and upper bounds. Therefore, all parameters that are assumed
- 157 to carry uncertainty are assigned to two different kinds of input distributions. We assume
- 158 distribution types found for the Meuse/Haute Marne URL site (France; [40]), and compare

159 the results to data gathered based on uniform priors for all parameters. As we are dealing
160 mostly with parameter uncertainties of the epistemic type, we are seeking to identify a method
161 that is robust in terms of prior parameter distribution assumptions.

- 162 2. In the second step, based on response quantities of interest, history match error metrics are
163 defined.
- 164 3. In the third step, we take the THM model and subject it to two different kinds of screening
165 designs (one-variable-at-a-time and Plackett-Burman) to obtain some initial understanding of
166 sensitivity.
- 167 4. Consequently, a proxy is built based on space-filling Latin-hypercube sampling. However,
168 before proceeding to real applications, we need to evaluate the size of the sampling space
169 required and estimate the resulting computational effort.
- 170 5. After checking essential quality measures, we use the proxy to perform a Monte-Carlo sam-
171 pling on the entire uncertainty space.
- 172 6. During the history matching step, only results are selected that follow observed (in our case,
173 synthetic) experimental data.
- 174 7. In an attempt to account for proxy errors, the filtered parameter sets are used to conduct
175 real model runs, that should support the agreement between the experimental and modeling
176 results.
- 177 8. The history-matched parameter sets are used for forecasts. As the idea of our test case is
178 to mimic a real-world experiment, we change the source term by increasing the power of the
179 heat source for the forecast. Since the analytical solution is not capable of describing this
180 change, we used a corresponding verified numerical models to conduct the forecast.
- 181 9. The uncertainty quantification in terms of percentiles of predefined performance or observa-
182 tion measures can then be done based on a forecast.
- 183 10. Finally, we present results from a global sensitivity analysis in terms of Sobol's indices that
184 are based on the proxy estimate in order to identify the main factors contributing to the
185 observed uncertainty. This is done both for prior and posterior distributions of input factors.

186 3.2. Workflow implementation: synthetic experimental data

187 The synthetic experimental data was created using a random seed selecting a set of input
188 parameter values uniformly within the lower and upper bounds of each input parameter (Tab. 1)

189 by evaluating the analytical model (AM). The data was gathered from a data sheet describing
190 the characteristics of the Opalinus Clay at the Mt. Terri site in Switzerland [41]. The generated
191 time-dependent response curves were varied with additional white noise (see Fig. 3).

192 3.3. Workflow implementation: parameter input for analysis

193 For system analysis, the same data with their corresponding input distributions were used
194 (Tab. 1 and SI Tab. 1). The functional forms for the input parameter distributions were taken
195 from a description of the DECOVALEX 2019 Task E specifications [42, 40]. Additionally, we also
196 performed an analysis based on the same input limits but using only uniform input distributions
197 to study the significance of the functional form of the input distributions. The water-related pa-
198 rameters a_w , ρ_w , K_w , c_w , and μ are relatively precisely known for given temperature and pressure.
199 However, as they are assumed to be constant in the analytical model (AM), we treat them here as
200 independent random variables with min/max values according to the temperature interval spanned
201 by $T_{min} = 290$ K and $T_{max} = 380$ K. As a matter of fact, such a treatment expands the uncer-
202 tainty space beyond the necessary scope as the underlying relation is known to be deterministic.
203 Nevertheless, as we are performing also a sensitivity analysis, it gives us also an insight into
204 which parameters can be regarded as constant due to their low sensitivity and for which parame-
205 ters we have to consider their deterministic relationship. In this study, we focused on five response
206 quantities namely temperature, pressure, displacement, and radial as well as circumferential stress
207 evaluated at an arbitrarily chosen observation point ($P = (0.5 \text{ m}, 0.0)$) and compared to the cor-
208 responding (synthetic) experimental data obtained as time series at the same location. In this
209 publication, we subject our investigation to a single observation point only, because of spherical
210 symmetry (i.e. the problem has only one effective spatial dimension) and the fact that both, the
211 synthetic experimental data as well as the history match-model (AM/NM) are of the same origin,
212 i.e., the information obtained at a single location should be enough to perform a sufficient history
213 match.

214 3.4. Workflow implementation: Design of Experiments (DoE)

215 In the applied workflow, methods of experimental design (DoE) are used to reduce the num-
216 ber of degrees of freedom in order to build a proxy model (PM) that can be used efficiently for
217 Monte-Carlo sampling and history matching. The degree of a possible agreement between experi-
218 mental and modeling (time-series) data permits the quantification of model uncertainties, whereas

219 the probabilistic analysis of the forecasted history matched model (*NM*) allows for uncertainty
 220 quantification of model and parameter uncertainties combined. The relative significance of dif-
 221 ferent input factors on the model output is investigated by means of a global sensitivity analysis
 222 based on the proxy model (*PM*) with particular reference to the history match error. The *PM* was
 223 built based on the space-filling Latin-Hypercube design. For the purpose of an initial parameter
 224 screening, the parameter space was sampled using a (folded) Plackett-Burman design [43] and a
 225 One-Variable-at-A-Time (OVAT) design at the domain bounds and around the 'best' values. This
 226 procedure is often referred to as local sensitivity analysis.

227 3.5. Workflow implementation: used software and libraries

228 The entire workflow was implemented in Python using the pyDOE2² library for the experimen-
 229 tal designs, GPy [44] for the proxy modeling of the history match error and SALib³ for the global
 230 sensitivity analysis (GSA). The entire workflow is wrapped around the multiphysics simulator
 231 OpenGeoSys⁴ [24, 45] and thereby ready for use with configurations of much greater complexity.

232 3.6. Problem Framing

233 The system under scrutiny is a greatly simplified model of a single canister of radioactive
 234 waste described by a point heat source in an infinite homogeneous isotropic porous fluid-saturated
 235 medium, as described in Section 2. In this case study, we analyzed the scalar quantities T and
 236 p as well as the u_r component of the displacement vector and the σ_{rr} and $\sigma_{\varphi\varphi}$ components of
 237 the stress tensor. As stated earlier, the response is only measured at one location, so we decided,
 238 therefore, to include both stress components as derived quantities in order to make the model
 239 sufficiently complex. The corresponding history match error metrics are then defined for each
 240 response quantity by

$$e^{\text{HM}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (d_i^{\text{obs}} - d_i^{\text{sim}})^2}. \quad (12)$$

241
 242 Here, n is the number of data/observation points in time, but also space, meaning that one
 243 error metric can, in principle, also be comprised of several observation locations. Under certain

²<https://pypi.org/project/pyDOE2/>

³<https://salib.readthedocs.io/>

⁴<https://www.opengeosys.org/>

244 (experimental) conditions, it might be useful to introduce an additional weighting factor for the
 245 data points. As we have only one observation point in this study, the sum is over a number of 2000
 246 time steps ($\Delta t = 5000$ s) used here for the synthetic data generation as well as for (NM-)modeling.
 247 d_i^{obs} corresponds to an observed datum, whereas the d_i^{sim} corresponds to the simulated data. In
 248 this study, d_i^{obs} are synthetic experimental data, as described in the previous section.

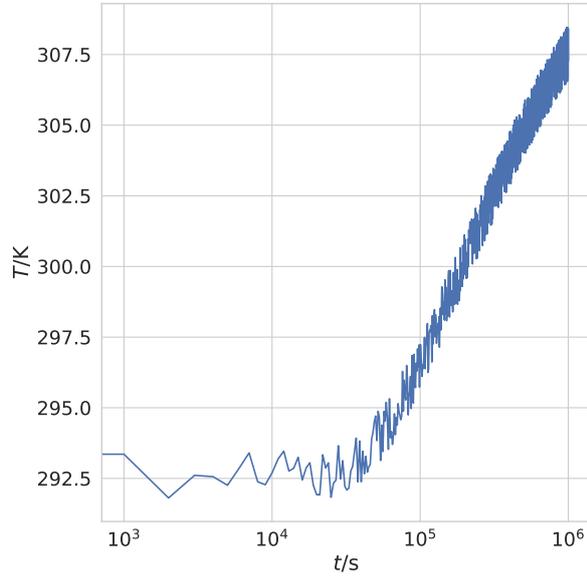
249 4. Results

250 4.1. Sensitivity Screening

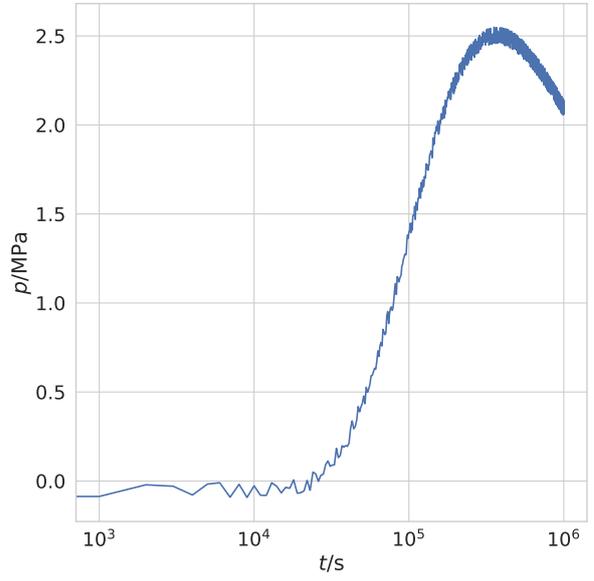
251 A local sensitivity analysis is conducted to get a first idea about the model (AM), i.e., its overall
 252 variability, its heavy hitters, and also any non-significant input parameters to be omitted during
 253 proxy building. We strive for a smaller number of inputs as they should result in better proxy
 254 quality if keeping the number of overall training samples fixed. Alternatively, seen from another
 255 perspective, we might be able to reduce the number of training samples resulting in a comparable
 256 proxy quality when having fewer inputs, i.e., a smaller uncertainty space. To serve this purpose,
 257 the size of the screening design should be taken much smaller compared to the runs needed for
 258 proxy training. More precisely, the required effort to neglect inputs should relate to a probable
 259 gain in time or proxy quality.

260 In this study, we start with a One-Variable-at-A-Time (OVAT) design to get a first insight and
 261 to create tornado plots (i.e., horizontal bar charts). The OVAT design was conducted in two ways.
 262 Both commence from a run with all parameters set to their baseline (referred to as 'best' in Tab. 1).
 263 The first type is intended to cover the entire parameter range, so changes in the response variables
 264 (e^{HM}) were obtained by changing the input parameters one at a time to their defined extreme
 265 values ('low'/'high' values in Tab. 1). In the second type of OVAT design, the changes around the
 266 baseline were reduced to a fraction of 100 and re-scaled afterward, i.e., the response was multiplied
 267 again by a factor of 100 to be comparable with the former design. However, the result of the latter
 268 type is more akin to a local sensitivity (tangent in the baseline), while the former evaluates the
 269 model (AM) at its actual extremes (secant). Additionally, we introduced a dummy parameter that
 270 is defined in the interval $[-1, 1]$ to obtain a clearer picture for distinguishing between significant
 271 and insignificant effects, which becomes important when applying statistical significance tests.

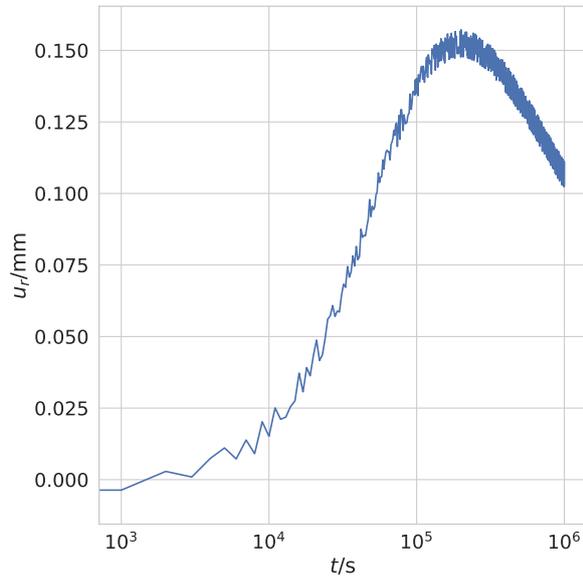
272 In Fig. 4 we show the tornado plots of the aforementioned types of OVAT designs for the
 273 pressure-related error metric e_p^{HM} . The tornado plots of the remaining response variables are



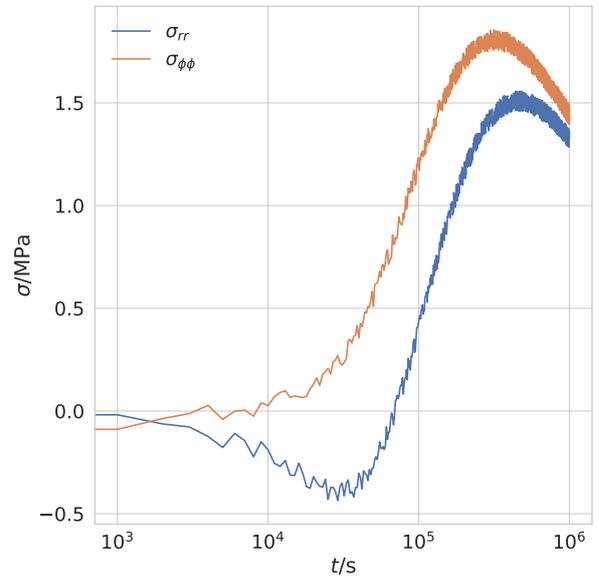
(a) Temperature



(b) Porepressure

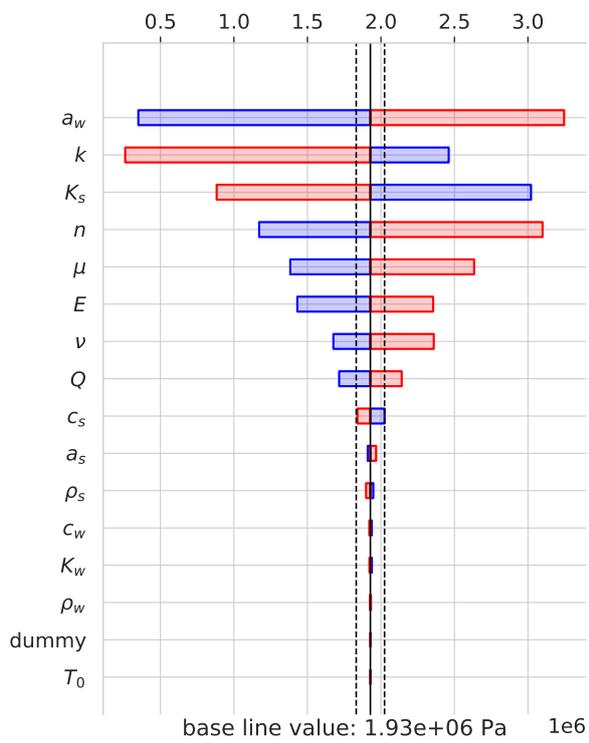


(c) Displacement

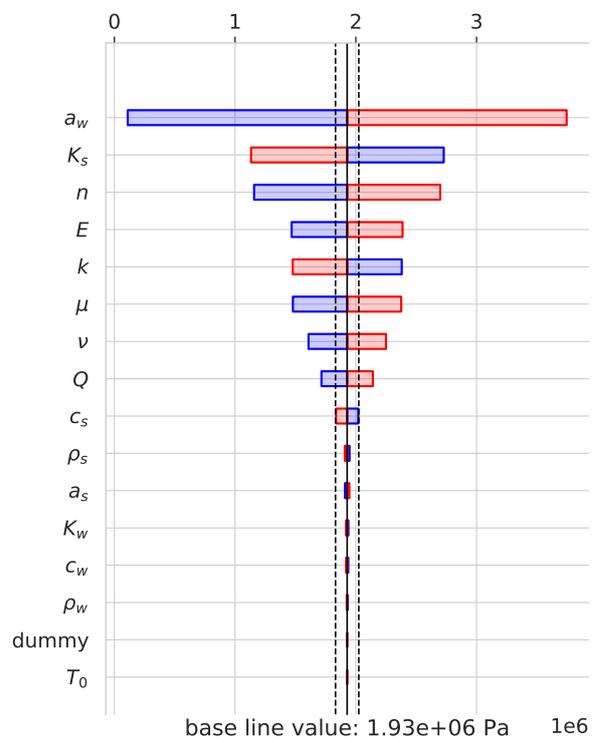


(d) Stress

Figure 3: Synthetic experimental response curves over the entire history match time span of 10^6 s (about 11 d) at Point $P = (0.5 \text{ m}, 0.0)$.



(a) bounds screening



(b) narrow band screening

Figure 4: Tornado plots of the pore pressure history match error (Pa) at point $P = (0.5 \text{ m}, 0.0, 0.0)$.

274 presented in Fig. 2-5 of the SI. Combining the information from all response variables, we can
 275 identify some parameters that have no or only a marginal impact (concerning the defined intervals)
 276 on the response variables. With a predefined triviality/significance margin of five percent, we find
 277 that ρ_w , ρ_s , c_w , K_w , and T_0 have only a small impact on all response variables, i.e., can be neglected
 278 later for proxy building.

279 The symmetry of the narrow band screening is due to the local perturbation of the base value
 280 and the subsequent re-scaling. Assymmetric effects due to a non-linear influence of model parame-
 281 ters or due to assymmetric 'high'/'low' values are only visible in the bounds screening variant of the
 282 tornado plots (Fig. 4). We also see that shifts in ranking occur due to the strong impact of the
 283 actual parameter variability on its associated sensitivity, an effect that is not captured by a local
 284 tangent.

It is essential to mention that the presented OVAT design has several drawbacks. One is that
 the system is probed only locally at a fixed position around the intermediate values. The second
 is that we do not test interactions between variables. To obtain more accurate results, we applied
 a folded Plackett-Burman design to screen the main effects as suggested in [38]. The sensitivity
 screening is then done employing a t-test on the regression coefficients using linear regression of a
 linear model (*LRM*)⁵:

$$e^{\text{HM}} = \beta_0 + \sum_{i=1}^k \beta_i x_i + \epsilon. \quad (13)$$

285 If studying also two-way (or higher) interactions, different designs need to be employed. Suitable
 286 experimental designs might be D-optimal and fractional-factorial designs keeping in mind that the
 287 minimum number of runs is given by the number of unknown coefficients of the regression equation.

288

289 While the Pareto chart in Fig. 5 shows a very similar behavior qualitatively as the tornado plot
 290 (Fig. 4), it appears as if only one parameter reached the critical t-value, i.e., can be considered as
 291 significant. However, a closer look at the coefficient of determination, and the F-statistics revealed
 292 that the sampling is not sufficient to exclude insignificant parameters because interaction terms
 293 were sampled but not included in the linear regression model (*LRM*).

294 As mentioned above, it is possible to switch to larger designs and also to account for interactions.

⁵For more on t-test and hypothesis testing in the context of regression analysis cf. *Regression Analysis by Example* [46].

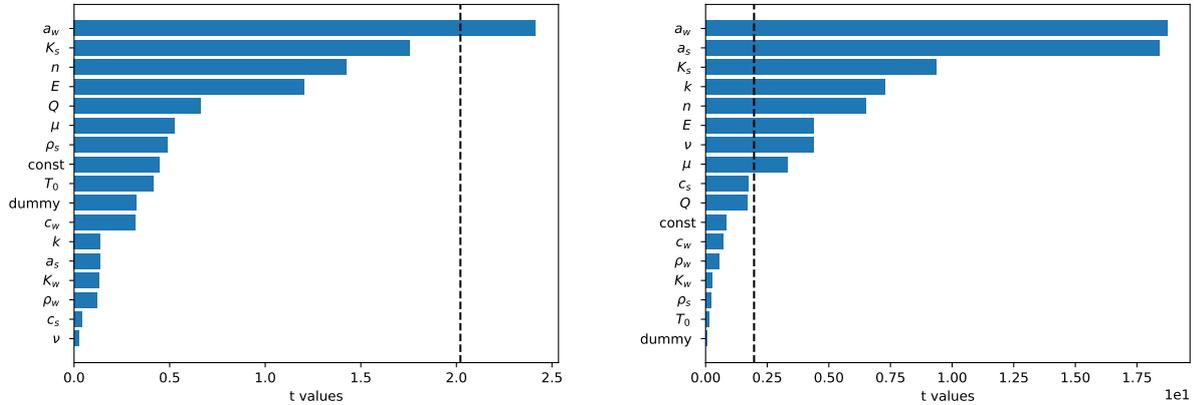


Figure 5: Pareto charts based on regression of a linear model for e_p^{HM} . Analysis of t values was conducted using a folded Plackett-Burman design (left) and latin-hypercube sampling (right) of the model domain (400 samples). The vertical line corresponds to the critical t-value for a significance level of $p = 0.05$.

295 As we are conducting this *initial* sensitivity screening for achieving gains in proxy performance,
 296 we require to stay well below one order of magnitude of the space-filling design we later use for
 297 proxy-building – otherwise the reduction of the uncertainty space would not pay off. However,
 298 we need to perform a space-filling design like Latin hypercube sampling (LHS) of the order of
 299 several hundred runs either way. Therefore, we can use the conducted design for proxy-building to
 300 perform an additional sensitivity screening in order to confirm our results so far. Using an ordinary
 301 LHS design with a size of 400 sampling points, we applied the t-test as presented above using a
 302 linear model (*LRM*) accounting only for main effects. Contrary to Plackett-Burman screening, the
 303 F-statistics support the claim of a significant influence of the linear fitting in general. Very much
 304 in agreement with the results from the OVAT screening, we find that ρ_w , ρ_s , K_w and c_w to have t-
 305 values below t_{crit} ($p \leq 0.05$) for all response error metrics, i.e. can be regarded as non-influential for
 306 all response parameters based on their given bounds of variation. When comparing both diagrams
 307 (Fig. 4 and Fig. 5, respectively), we find that the order of sensitivity of most parameters changes,
 308 which comes from different contributions of certain parameters at different locations when using
 309 a space-filling design. Another point, we want to stress here, is that in all cases, the dummy
 310 parameter ranks above other parameters while staying well below the critical t-value, which gives
 311 us confidence that those parameters do not have a significant impact on the model *AM*. The fact
 312 that the t-value is not precisely zero is not a surprise as we are trying to fit a linear model (*LRM*)
 313 with a manageable amount of samples. Even though using a space-filling design, we might overlook

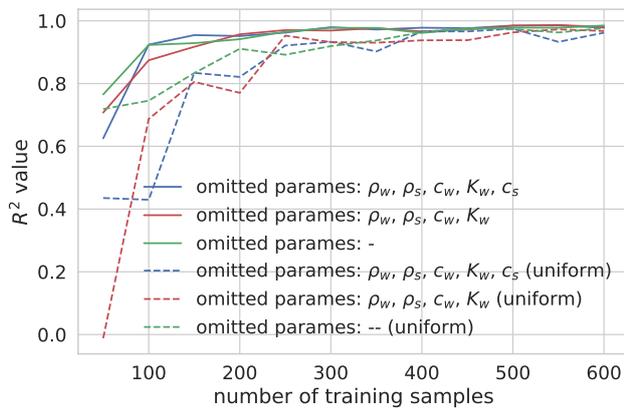
314 some effects that come from, e.g., a smiley face pattern or effects that stem from interactions of
315 input parameters. However, at this stage only a preliminary screening is intended and a more
316 thorough sensitivity analysis will be conducted later using the proxy (PM) and sampling-based
317 methods that also provide much better coverage of the uncertainty space [29, 12] and that allow a
318 better quantification of relative effects.

319 4.2. Proxy Building

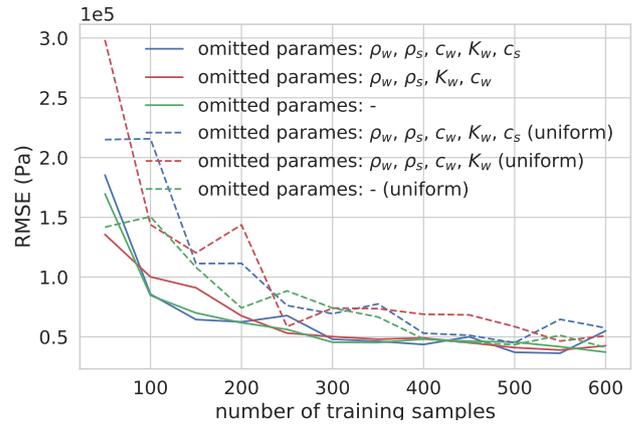
320 For proxy building, as stated in the previous paragraph, a space-filling design like Latin hyper-
321 cube sampling (LHS) needs to be employed. As proxies, typical choices are polynomial, splines,
322 kriging, or neural networks. A very suitable choice for our purpose is Gaussian process regression
323 (often referred to as kriging in the geosciences) as every training point can be recovered precisely
324 by the proxy, so it is thereby expected to provide better accuracy compared to parametric proxies.
325 Comparisons for different types of proxies can be found elsewhere [38, 47]. At this stage, we take
326 all parameters with assigned uncertainty with us, as we want to study what kind of influence the
327 omission of some of them will have on the proxy quality. The proxy was built based on 50, 100,
328 150, ..., 400 training samples plus 200 testing samples of a space-filling Latin-hypercube design
329 to observe the convergence behavior of the coefficient of determination in order to determine a
330 required number of training samples.

331 In Fig. 6(c), we plotted the time needed to build and apply the proxy to the testing samples,
332 together with the R^2 (Fig. 6(a)) and RMSE quality measures (Fig. 6(b)) for e_p^{HM} over a varying
333 number of training samples. Different curves correspond to different classes of a priori input dis-
334 tributions (non-uniform vs. uniform assumptions) and different sets of omitted input parameters.
335 What we see at first glance for the proxy quality measures, is that the R^2 /RMSE values improve
336 quite strongly between 50 and 200 training samples, while changes tend to converge for higher
337 numbers. Keeping in mind that for each training and testing sample, a full run of the typically
338 very costly NM -model is required, the search for a sample number compromising between accu-
339 racy and computational effort becomes obvious. Therefore, one would prefer the application of
340 a sequential sampling strategy [48] and stop increasing the number of training samples when a
341 justifiable R^2 /RMSE value is reached.

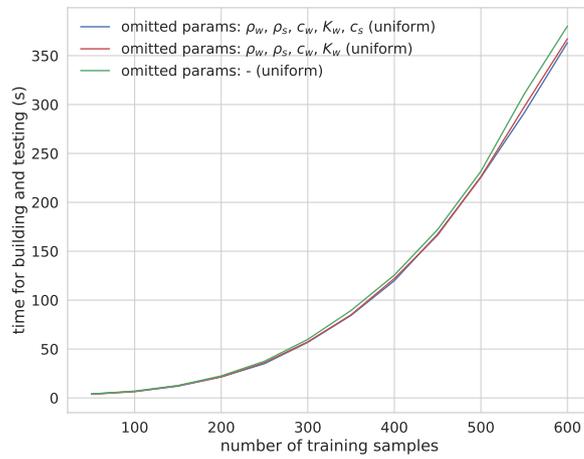
342 As can also be seen from the plots, we can clearly distinguish between uniform and non-uniform
343 parameter distributions, although the differences vanish for a higher number of sampling points.



(a) Coefficient of determination (p)



(b) Root-mean-square-error (p)



(c) Time for proxy building.

Figure 6: Performance of the e_p^{HM} kriging proxy for a varying number of input parameters and training samples.

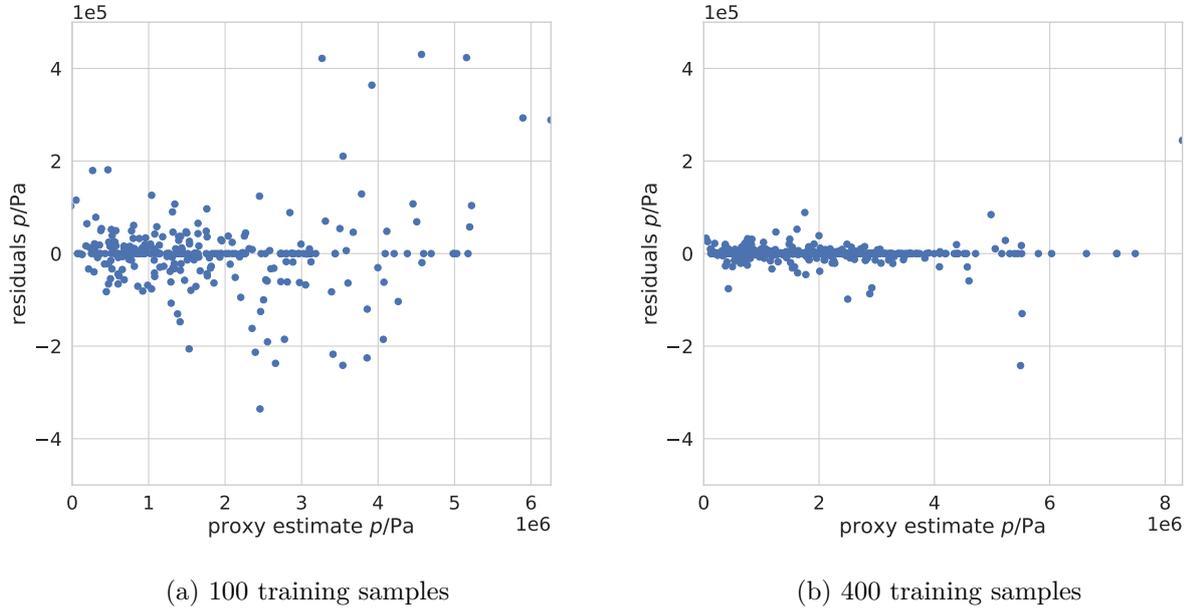


Figure 7: Residuals of the e_T^{HM} error metric based on 100 and 400 training samples.

344 This means, with respect to proxy quality, that the drawback of not knowing the exact distribution
 345 can be compensated by choosing the size of the Latin-hypercube design big enough. On the other
 346 hand, no clear trend is visible for the omission of non-influential parameters. Indeed, the difference
 347 between uniform and non-uniform parameter distributions is somewhat expected, as in the uniform
 348 case, a greater parameter space is covered, leading to a lower sampling density on the average.
 349 The not-shown response quantities also confirm these trends (see Fig. 8-9 SI for further details).
 350 Similar to the quality measures, the time required for proxy building and testing is also not very
 351 sensitive to parameter omission. As R^2 and RMSE give us a rough sense of the proxy quality, we
 352 also had a look at how the residuals are, in fact, distributed (Fig. 7). The residuals are plotted
 353 versus their associated proxy estimates for 100 and 400 training samples. Both plots show only
 354 slight heteroscedastic behavior, mainly an increasing variance for a higher proxy estimate. While
 355 greater heteroscedasticity could pose a real problem, slight heteroscedasticity is often unavoidable.
 356 Another important measure of proxy quality is whether the proxy preserves essential mathematical
 357 and physical properties. One such property of the history match error is that it is strictly positive.
 358 In Fig. 8 one sees that a small fraction, especially for smaller sample sizes, does not fulfill this
 359 criterion. However, for the σ_{rr} - proxy built using 400 training samples only a fraction of $< 10^{-3}$

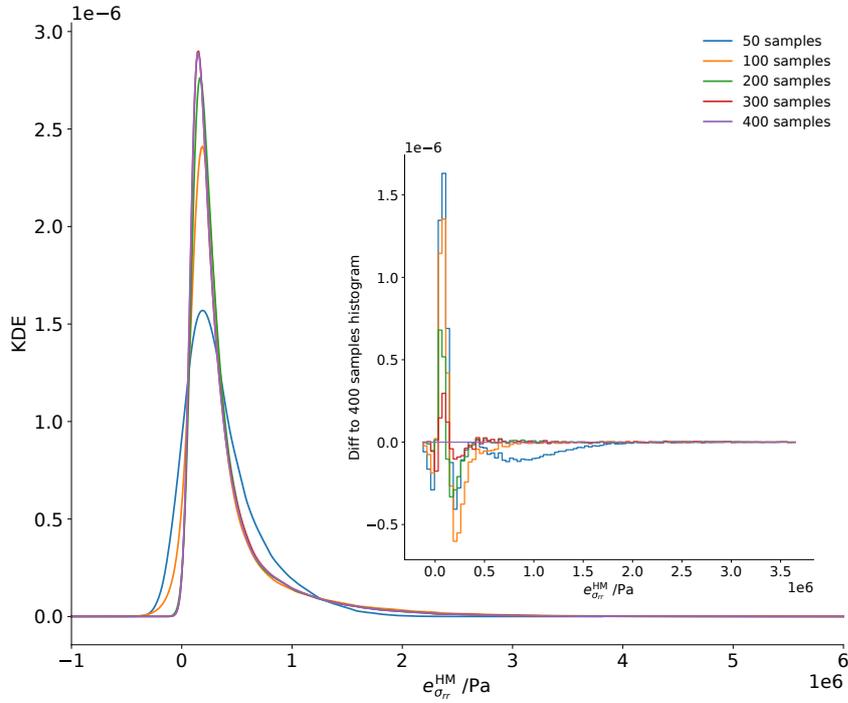


Figure 8: Kernel density estimate of the $e_{\sigma_{rr}}^{\text{HM}}$ proxy error metric and their histogram of differences to 400 samples reference based on 200,000 MC samples.

360 is negative, keeping in mind that the expectation value of $e_{\sigma_{rr}}^{\text{HM}}$ is $\approx 10^5$ Pa. To conclude our
 361 analysis, a sample number of 400 without omitting any input seems to be a reasonable choice for
 362 our further investigation as no significant improvement could be achieved by either increasing the
 363 sample numbers or omitting sets of marginally influential parameters.

364 4.3. Monte-Carlo Sampling and history match Filtering

365 Based on the choice of the previous paragraph, the proxy model (*PM*) is ready for Monte-Carlo
 366 sampling. 200,000 Monte-Carlo samples are drawn from the prior 'non-uniform' and 'uniform'
 367 distributions and evaluated for each proxy error metric.

368 For testing purposes, we applied three increasingly strict filter criteria (Tab. 2) to the sampling
 369 results and analyzed the output (Fig. 9) as well as the corresponding history-matched input param-
 370 eter distributions (Figs. 10 and 11) for non-uniform and all uniform priors. Thereby, we selected
 371 only samples that satisfied the history match criteria for all variables simultaneously. Whereas
 372 blue corresponds to the prior distributions, orange, green, and red denote the different conditions
 373 as defined in Tab 2. All histograms are re-scaled after filtering in order to be visible in the figures.

374 Taking first a look at Fig. 9, we see that except for the temperature, all history match curves
 375 are lying between values around zero and the cut-off condition smaller than the prior histogram
 376 reaching its maximum. The temperature-'anomaly' can be explained by the fact that the range
 377 of possible solutions is comparable to the noise of the synthetic experimental curve (cf. Fig. 3).

378 Therefore, the history match error is very much affected by the noise, which also leads to the shift
379 of around 0.5 K in the histogram. Comparing the results to the data obtained from uniform prior
380 distributions, we see that the posterior distributions seem to behave very similarly. In other words,
381 if the prior uncertainty space is much bigger than the space obtained by the filter, then the exact
382 forms of the prior parameter distributions become irrelevant when sampled appropriately. Another
383 more subtle point is that the condition to satisfy all criteria simultaneously leads to a significant
384 reduction of the number of matched curves compared to the histogram’s area that would remain
385 after application of the individual quantity’s criterion alone.

variable		test cond. 1	test cond. 2	test cond. 3	final cond.	RMSE	unit
e_T^{HM}	<	3	2.5	2.0	2.0	$1.5 \cdot 10^{-3}$	K
e_p^{HM}	<	$1.5 \cdot 10^5$	$1.0 \cdot 10^5$	$0.5 \cdot 10^5$	$0.3 \cdot 10^5$	$0.5 \cdot 10^5$	Pa
$e_{u_r}^{\text{HM}}$	<	$1.5 \cdot 10^{-5}$	$1.0 \cdot 10^{-5}$	$0.5 \cdot 10^{-5}$	$0.3 \cdot 10^{-5}$	$0.5 \cdot 10^{-5}$	m
$e_{\sigma_{rr}}^{\text{HM}}$	<	$1.5 \cdot 10^5$	$1.0 \cdot 10^5$	$0.5 \cdot 10^5$	$0.3 \cdot 10^5$	$0.35 \cdot 10^5$	Pa
$e_{\sigma_{\varphi\varphi}}^{\text{HM}}$	<	$1.5 \cdot 10^5$	$1.0 \cdot 10^5$	$0.5 \cdot 10^5$	$0.3 \cdot 10^5$	$0.35 \cdot 10^5$	Pa

Table 2: History match filtering conditions for all response parameters. The testing conditions were used to investigate the effect of varying filter sizes using the *AM*-model. The final condition together with the *PM*-RMSE is used to demonstrate the workflow functionality including model runs for history match and forecast(*AM/NM*).

386 In a second step, we analyzed the input parameter combinations that were used to generate the
387 response proxies that survived the history match filtering (‘filtered priors’). For each parameter
388 and filter condition, the corresponding distributions are depicted in Figs. 10 and 11. Here, we see
389 that for some parameters like E , the filter does not seem to have any impact on the distributions,
390 whereas for others (like K_s , a_s , k or μ) the filter seems to restrict the domain of definition. Indeed,
391 the latter parameters were shown to be heavy hitters in sensitivity screening. For these parameters,
392 we also see a better agreement between the input used to generate the experimental data and the
393 posterior distribution with the tightest criterion. Parameters that retain their distributions after
394 filtering were found to be less influential, which is very much in agreement with what one would
395 assume. It is also worth noting that the posterior distributions found for all uniform priors agree
396 very well with the ones for non-uniform priors for the influential parameters. This gives rise to
397 the conclusion that the exact form of the prior distributions is not very relevant for the history
398 match results. Another thing that attracts our attention is that, for k and μ , we find a kind of
399 multimodal behavior after history matching that we see, especially for non-uniform priors. This
400 behavior, as well as the fact that the sampling maxima of the tightest match criteria do not

401 coincide with the input of the synthetic model (*AM*) run, we attribute to parameter interactions
 402 that cannot be disregarded. In general, the history matching workflow is also quite suitable to serve
 403 as an alternative to optimization algorithms for the purpose of parameter estimation. However,
 404 the latter discussed behavior also makes clear that the reverse problem of parameter estimation
 405 of such a highly non-linear model is ill-posed. For further analysis of parameter interactions, we
 406 allude to the global sensitivity analysis.

The impact of different prior distributions as well as the effect of different filter sizes on the direct values of the response functions are also investigated using the *AM*-model. For this analysis, we used the last time step of the history match and calculated the corresponding CDFs based on the testing filters defined in the previous section. The cumulative distribution functions (CDF) were evaluated at point $P = (0.5 \text{ m}, 0.0, 0.0)$ and time $t = 5 \cdot 10^6 \text{ s}$. In Fig. 12, the CDFs were presented based on prior distributions before and after applying the three different filter conditions. For all three conditions and all response functions, we find that the co-domain is significantly reduced, and we are able to give percentiles for each quantity from which we can select representative models. One major property that such a workflow needs to satisfy in order to demonstrate its robustness is that its results should not depend too much on arbitrarily chosen values. While it is quite obvious, that the filter size might have a greater impact on p10 and p90 percentiles, we assume, that this influence tends to be rather small for the p50 percentile. To assess the validity of this assumption, we calculated the relative change of the p50 value for each filter with respect to no filter as follows:

$$p^X = \frac{|X_{p50}^{\text{filter}} - X_{p50}^{\text{no filter}}|}{X_{p50}^{\text{no filter}}} \quad (14)$$

407 The corresponding values are given in Tab. 3. Here, we see that for non-uniform priors, the p50
 408 value varies up to six percent, while the relative changes for all uniform priors are all systematically
 409 greater and amount up to 9 percent for the displacement. However, it is not obvious that the
 410 expectation for all uniform priors is subject to greater fluctuations in all cases, nor does our
 411 example prove that we get more reliable results for non-uniform priors. While these changes can
 412 be regarded here as quite small, they might matter in some cases and remind us that the choice of
 413 the filter size is somehow subjective, and it necessitates a more thorough analysis when discussing
 414 concrete safety functions.

415 After having obtained an overview, we are able to define the actual history match criteria
 416 with which we will perform real model runs (*AM/NM*) and consequently, the forecast. To define

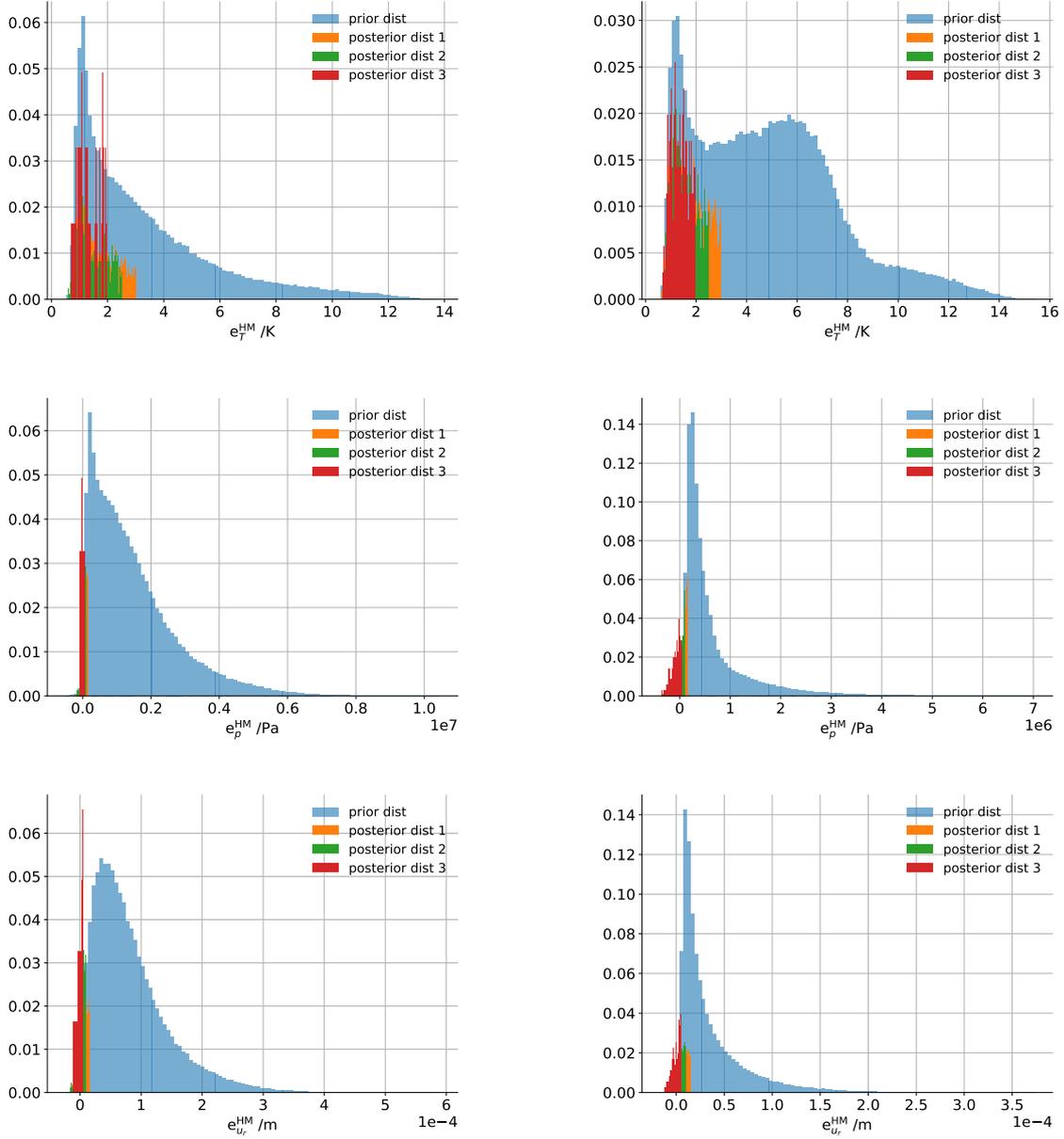


Figure 9: Distributions of the history match error responses for non-uniform (left) and all uniform priors (right) and the filter conditions listed in Tab. 2.

417 these conditions, one has to take also into consideration that several error sources should also
 418 have an impact on the defined confidence interval. Most important (i) is the (*AM/NM*) model
 419 error (e.g., heterogeneities or other processes that have to be taken into account), (ii) followed
 420 by the proxy error and (iii) last but not least, the sampling error. As the model error cannot be
 421 precisely determined, the choice of the history match thresholds coming from this contribution is

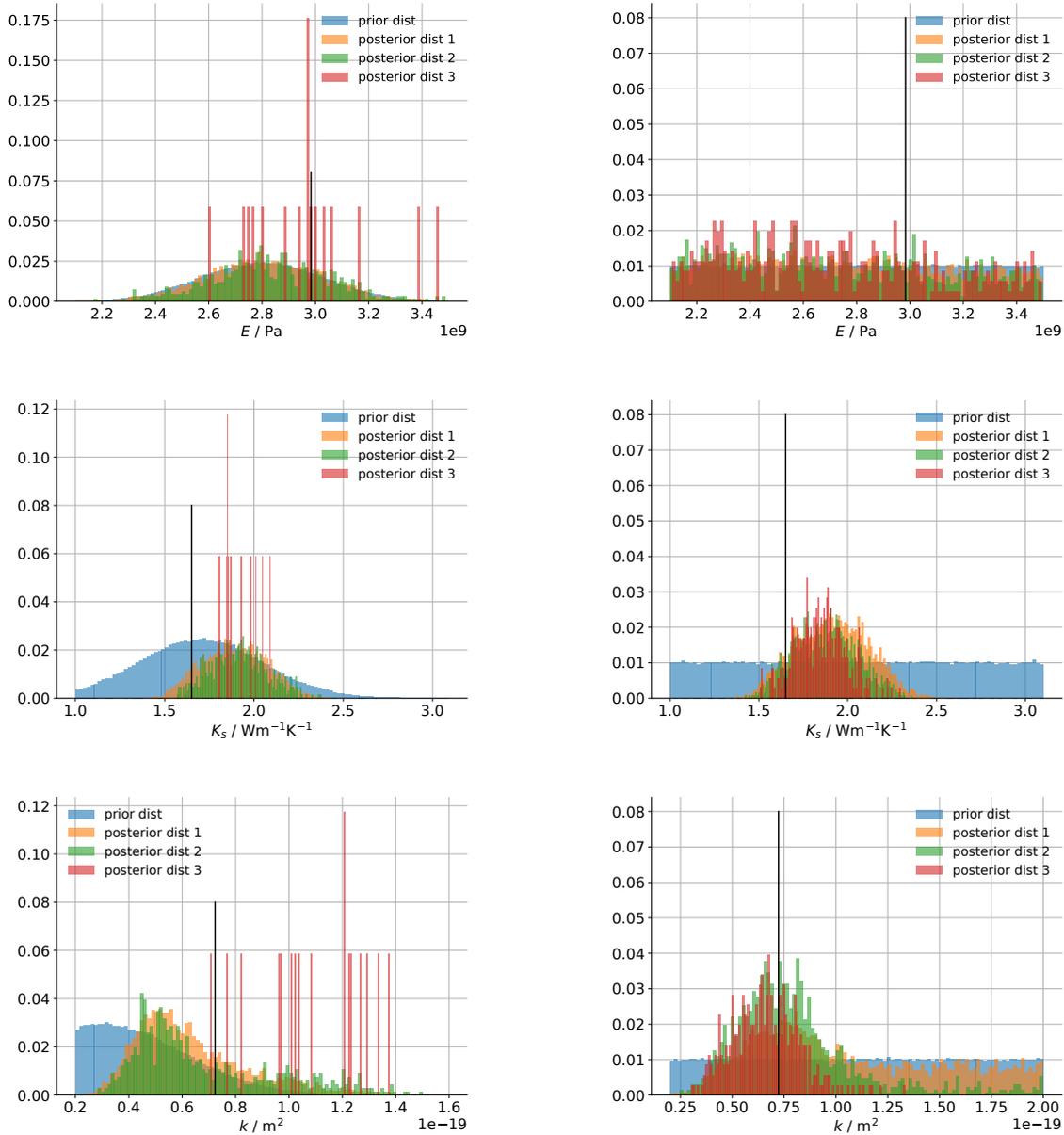


Figure 10: Parameter estimation for a selection of input quantities from Monte-Carlo sampling and after history match filtering for non-uniform (left) and all uniform (right) priors with the filter conditions listed in Tab. 2. The black line marks the input parameter for the creation of the synthetic experimental data.

422 somehow subjective. For our problem case, the criteria are given by the fifth column in Tab 2 (final
 423 conditions). Additionally, the proxy error was taken into account by incorporating the RMSE (sixth
 424 column, Tab 2) twice into the definition of the threshold [38]. One more technical criterion is that
 425 we need a sufficient number of surviving models after the intersection of all error metric thresholds

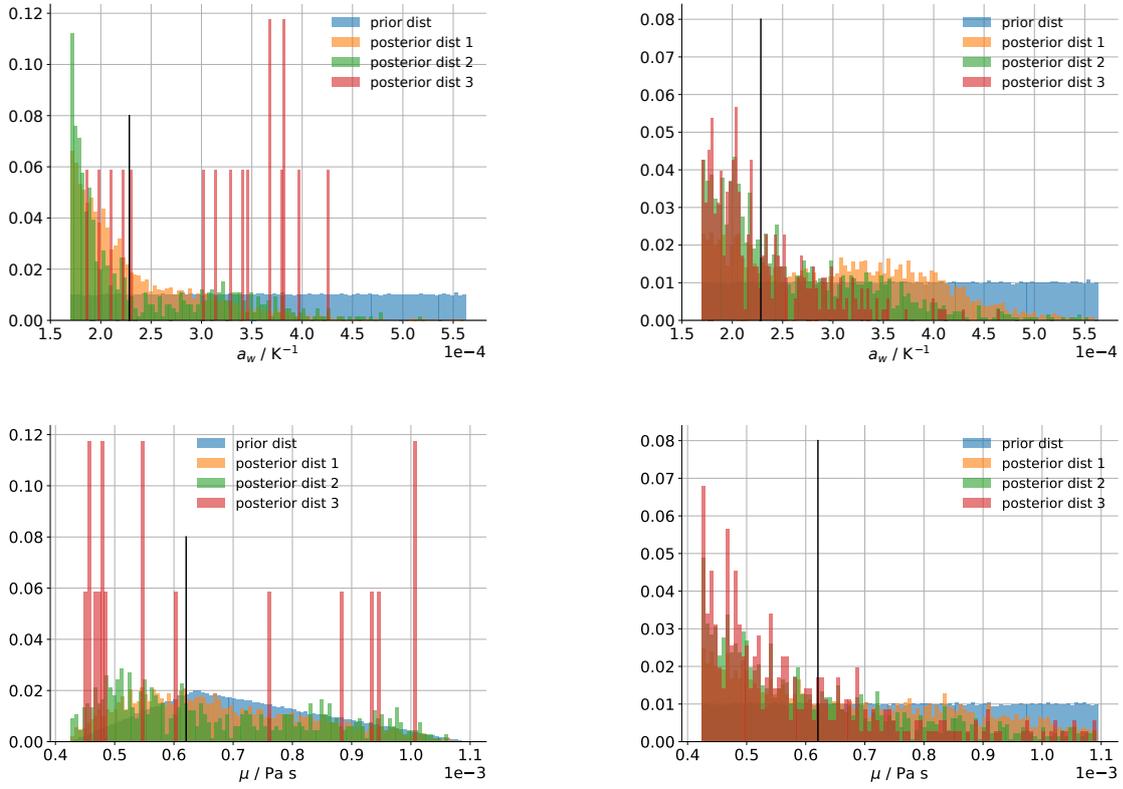


Figure 11: Parameter estimation for a selection of input quantities from Monte-Carlo sampling and after history match filtering for non-uniform (left) and all uniform (right) priors with the filter conditions listed in Tab. 2. The black line marks the input parameter for the creation of the synthetic experimental data.

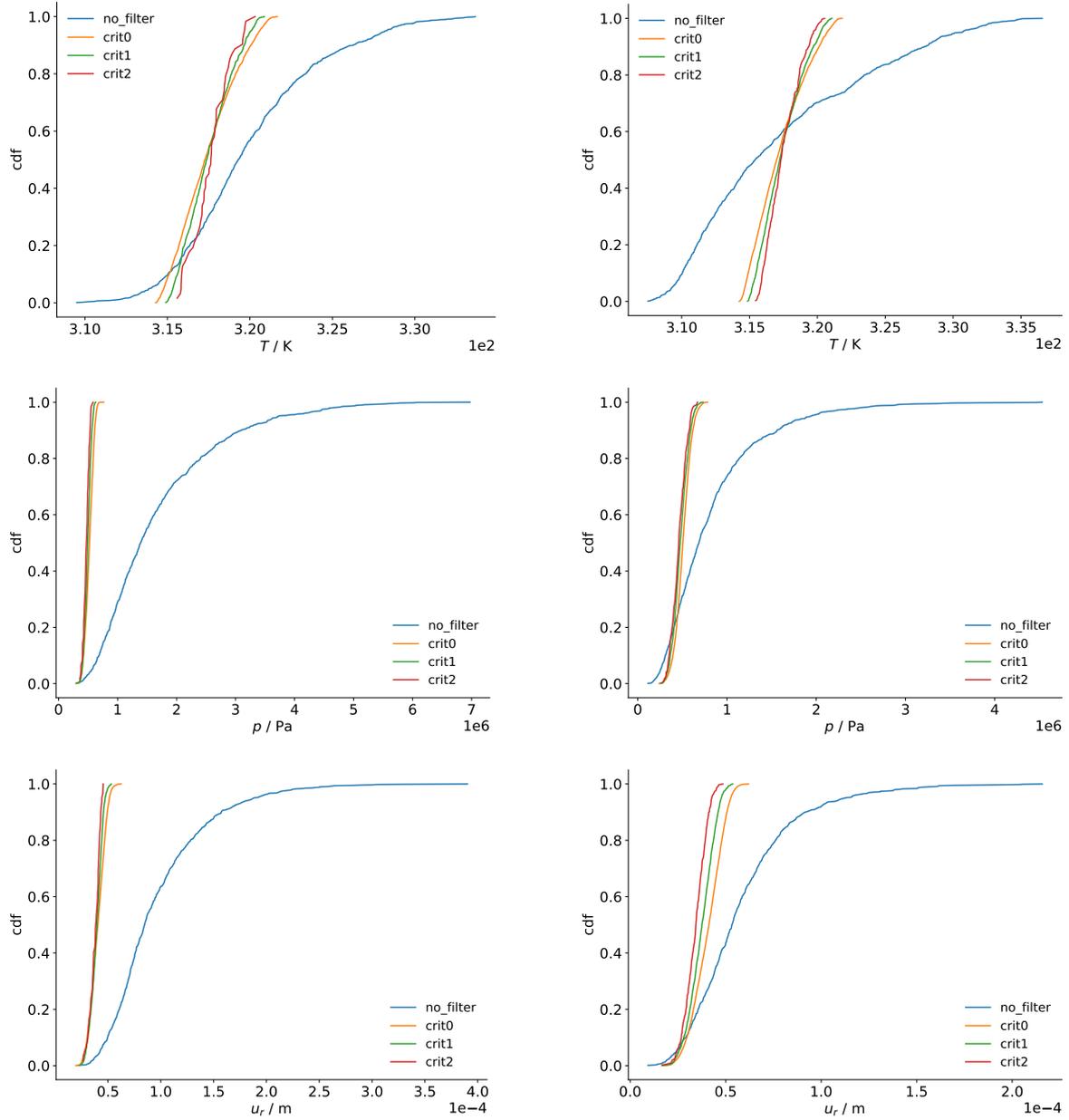


Figure 12: Cumulative distribution functions of the last history match time step obtained from Monte-Carlo sampling and after applying the testing history match conditions, see Tab. 2.

Table 3: Relative p50 values for all three testing filters with respect to the p50 value obtained from Monte-Carlo sampling.

	p^T		p^p		p^{u_x}	
	non-uniform	uniform	non-uniform	uniform	non-uniform	uniform
crit0	0.006866	0.003734	0.632036	0.266558	0.514512	0.251206
crit1	0.006673	0.004797	0.641845	0.286974	0.518569	0.315963
crit2	0.007089	0.006498	0.655329	0.319417	0.577460	0.340364

426 to perform a statistical analysis afterwards if probabilities are to be attached to the outcome. With
427 the matched criteria, 'full' model (AM) runs were conducted again, to check whether the proxy
428 estimates and history match criteria are in sufficient agreement with our assumptions and our
429 conception of the history match (Fig. 13).

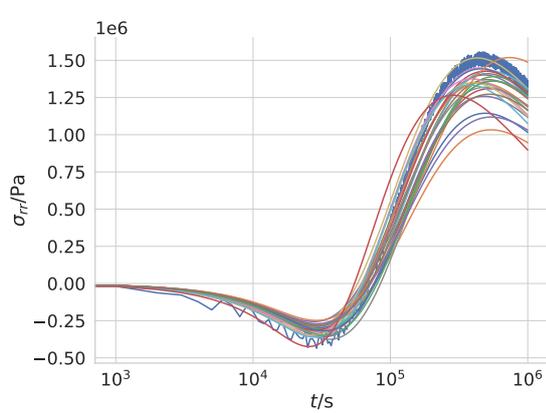
430 4.4. Forecast

431 Analogous to the (AM)-model runs of the history match, we repeated the same calculation
432 this time with an equivalent numerical model (NM) with a subsequent forecast time ($1 \cdot 10^6$ s –
433 $-5 \cdot 10^6$ s). From $1 \cdot 10^6$ s to $1.1 \cdot 10^6$ s the power of the heat source is linearly increased from
434 300 W to 600 W and kept constant until the end of the run. The corresponding stress response
435 over the entire simulation time is presented in Fig. 13. The changes due to the altered source
436 term conditions are clearly visible and can be confirmed or rejected *a posteriori* by experiments
437 or monitoring data. In case of a rejection, the model would require further adaptations. In this
438 study, we assume for simplicity that the chosen observation quantities are somewhat representative
439 of the safety functions to be monitored. The corresponding CDFs of the last forecast step, yield
440 analogous results as shown in Fig. 12 and are not subject to further analysis in this study.

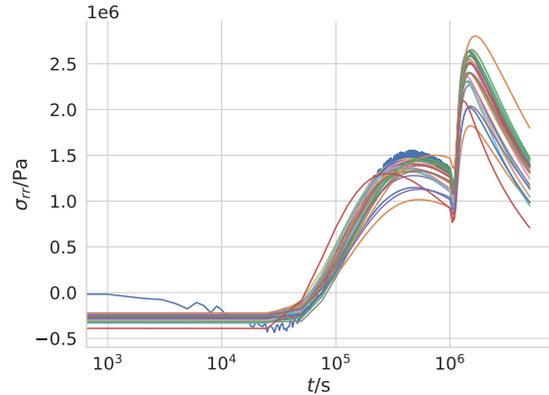
441 4.5. Proxy-Based Global Sensitivity Analysis

442 Now that model uncertainty has been investigated, it is of practical interest to attribute this
443 uncertainty to individual parameters or their combinations. In other words, one would like to
444 understand how variations of the model input affect the history match error of the response func-
445 tions. For that purpose, we analyzed the PM in terms of Sobol's indices by using their Monte-Carlo
446 estimates.

447 For this purpose, we sampled the input space between their *min* and *max* values using the
448 sampling scheme of Saltelli [29]. First and second-order indices were calculated using their Monte-



(a) history match (analytical model)



(b) history match and forecast obtained from numerical model

Figure 13: σ_{rr} as function of time at point $P = (0.5 \text{ m}, 0.0, 0.0)$. a) Curves from *AM*-model satisfying all history match criteria. b) σ_{rr} results obtained from numerical model (*NM*) for the same parameter set including forecast with modified source term. In both cases the experimental curve is given in blue in background.

449 Carlo estimates. For all response parameters, we used a sampling size of 10,000, which was shown
 450 to be sufficient for the width of the 0.95 confidence interval to be well below 0.025 for all indices
 451 (cf. Fig. 10 SI).

452 All in all, the global sensitivity analysis for our five error metric proxies provides a clear picture
 453 of the influence of single parameters on the model output. The results for both first- and second-
 454 order indices are given in Fig. 14. For the second-order indices, we present only values that exceed
 455 the error margin of 0.025. The bounds for the posterior analysis were estimated by the extreme
 456 values found during parameter estimation after applying the final filter condition (cf. Fig. 10
 457 and Fig. 11). One general trend that attracts our attention is that more parameters become
 458 influential after applying the filter condition. Whereas the thermal conductivity K_s overwhelmingly
 459 dominates the influence on the temperature, after filtering, also other factors become important
 460 on the temperature. The analysis of the main effects is very much in agreement with the findings
 461 during parameter screening (Sec. 4.1). There, we found that ρ_w , ρ_s , K_w , c_w are non-influential
 462 and can be neglected for uncertainty analysis. The analysis also confirms the findings of Sec. 4.3:
 463 heavy hitters like K_s , k , or a_w changed their behavior after filtering. Looking at the second-order
 464 indices also reveals that the shift for K_s and the 'multimodal' behavior, we found for k , μ or a_w
 465 can be very much attributed to interaction effects between different parameters. At this point,
 466 it is important to note again that due to the challenges the analytical solution poses in terms of

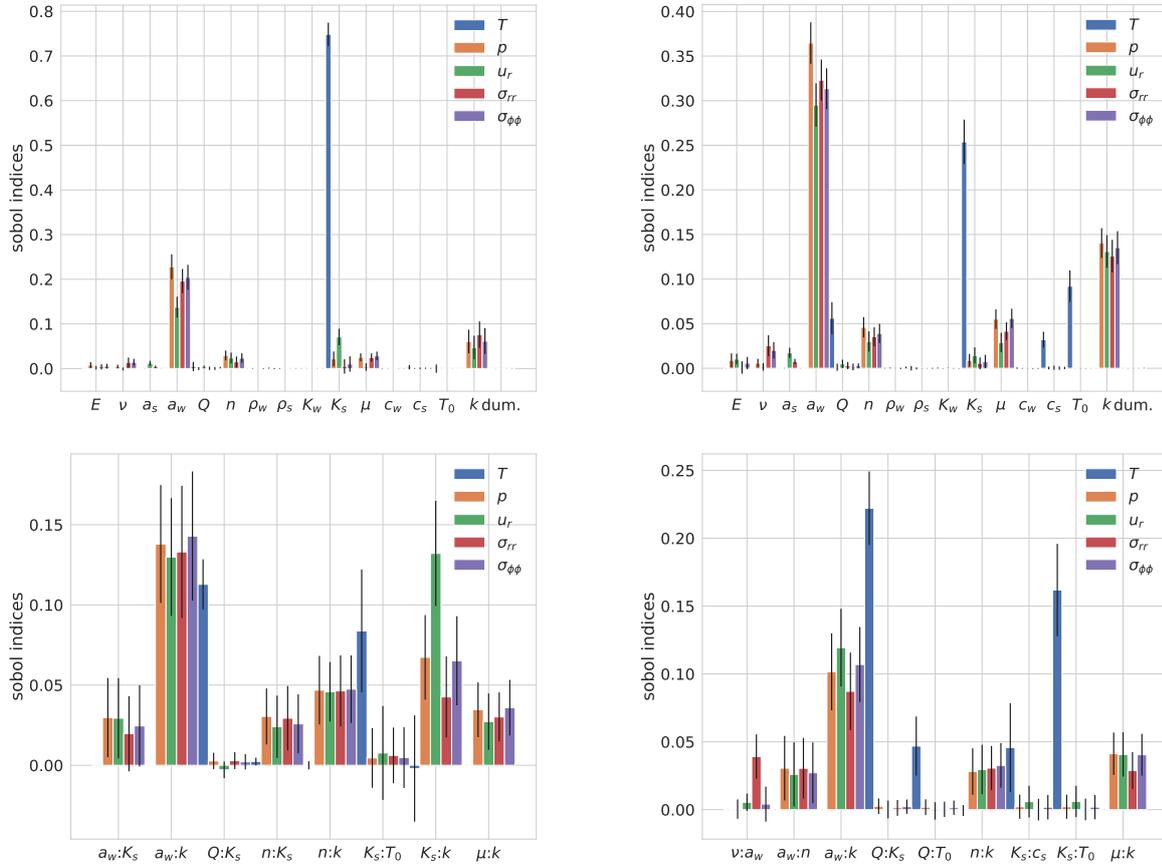


Figure 14: Sobol indices based on prior assumptions (left) and posterior (right) bounds. The upper graphs contain only indices of first order effects, while the lower contain second order interactions. For the second order terms, only combinations with effects greater than 0.025 are shown.

467 constant parameter input, we treated all of the water-related constants as uncertain themselves.
 468 From the GSA we can also conclude that we can treat ρ_w , K_w , and c_w as constants, while the
 469 functional dependence of μ and a_w seems to be relevant for further analysis utilizing finite element
 470 software.

471 5. Conclusions

472 In our work, we scrutinized the applicability of the Design of Experiments (DoE)-based assisted
 473 history matching for uncertainty quantification in linear coupled thermo-hydro-mechanical models.

474 In our manuscript, we used an analytical model of a simplified geotechnical problem in the form
 475 of a disposal cell containing heat-emitting radioactive waste emplaced in an isotropic fluid-saturated
 476 medium under realistic parameter conditions.

477 As these parameters often cannot be described in terms of a known probability distribution
478 function (PDF), we compared PDFs based on expert opinion with a case of all uniform input
479 parameter distributions. The most important findings of the presented work are:

- 480 • While in order to find a good history match, the filtered response must be covered by a
481 valid input parameter range, it was shown in the present study that the exact form of the
482 input parameter distribution becomes less critical. The result that the determining factors
483 are the filtering conditions instead of the input parameter distribution makes the approach
484 particularly interesting for our purpose in dealing with epistemic uncertainties as precise
485 distribution forms are not known for the entire system.
- 486 • It was found that a detailed exploration of input parameter distributions before modeling is
487 beneficial in reducing the uncertainty space and improving the proxy quality. At the same
488 time, precise knowledge of all parameters is not required in order to obtain good history
489 matched results. However, this comes at a cost in computing time as a broader sampling
490 might be required.
- 491 • One disadvantage of the workflow, at first sight, is the somehow subjective choice of the
492 history match filtering conditions. As we could show, the filter size can have an effect on
493 the stochastic outcome (like percentiles), which can be crucial. The filter size describes the
494 effect of uncertainty reduction due to the agreement of model output with experimental data.
495 Therefore, care should be taken in quantifying these discrepancies for defining the history
496 match thresholds. In case of doubt, one should assume them to be less tight. However, this
497 is also part of a more general problem when dealing with the quantification of uncertainties
498 that are of epistemic origin and involve model predictions: The meaning of exact numerical
499 values is often overrated, and the usage of rigorous mathematical concepts often obscures the
500 fact that the underlying problem eludes an exact quantification.
- 501 • As the history matching procedure reduces the uncertainty space significantly, it also affects
502 the sensitivity, i.e. the relative ranking of input parameters. Therefore, conclusions drawn
503 regarding sensitivity analyses prior to an experiment may have to be re-evaluated after (more)
504 experimental data has become available.

505 In our manuscript, we showed that the workflow is particularly suitable for uncertainty quantifi-

506 cation, sensitivity analysis, and model validation in geotechnical applications like radioactive waste
507 repositories. However, before turning directly to real-world applications, the conceptual validity
508 and computational feasibility of even more complex models incorporating non-linear phenomena
509 (e.g. equations of state, material behavior) and spatial heterogeneities need to be demonstrated.

510 Acknowledgments

511 This work was funded in parts by the German Federal Ministry of Education and Research
512 (BMBF) as part of the project Integrity of nuclear waste repository systems – Cross-scale system
513 understanding and analysis – iCross (02NUK053E) and Impulse and Networking Funds of the
514 Helmholtz Association (SO-093). We are also very grateful to the OpenGeoSys developer team for
515 their enthusiastic, continuous work on further developing and improving the OGS platform for the
516 scientific community.

517 References

- 518 [1] J. C. Helton, J. D. Johnson, C. J. Sallaberry, C. B. Storlie, Survey of sampling-based methods for uncertainty
519 and sensitivity analysis, *Reliab Eng Syst Saf* 91 (10-11) (2006) 1175–1209.
- 520 [2] C.-F. Tsang, O. Stephansson, J. Hudson, A discussion of thermo–hydro–mechanical (THM) processes associated
521 with nuclear waste repositories, *Int J Rock Mech Min Sci* 37 (1-2) (2000) 397–402.
- 522 [3] F. Bernier, F. Lemy, P. De Cannière, V. Detilleux, Implications of safety requirements for the treatment of
523 THMC processes in geological disposal systems for radioactive waste, *J Rock Mech Geotech Eng* 9 (3) (2017)
524 428–434.
- 525 [4] X.-L. Li, J. Lanru, P. Blaser, Impact of Thermo-Hydro-Mechanical- Chemical (THMC) processes on the
526 safety of underground radioactive waste repositories. Proceedings of the European Commission TIMODAZ-
527 THERESA International Conference, Luxembourg, 29 September - 1 October 2009, in: Proceedings of the
528 European Commission TIMODAZ-THERESA International Conference. Luxembourg: Publications Office, no.
529 October, 2009, pp. 489–494.
- 530 [5] C. F. Tsang, J. D. Barnichon, J. Birkholzer, X. L. Li, H. H. Liu, X. Sillen, Coupled thermo-hydro-mechanical
531 processes in the near field of a high-level radioactive waste repository in clay formations, *Int J Rock Mech Min*
532 *Sci* 49 (2012) 31–44.
- 533 [6] J. C. Helton, Uncertainty and sensitivity analysis techniques for use in performance assessment for radioactive
534 waste disposal, *Reliab Eng Syst Saf* 42 (2-3) (1993) 327–367.
- 535 [7] R. G. Ghanem, P. D. Spanos, *Stochastic Finite Elements: A Spectral Approach*, 1991.
- 536 [8] G. Stefanou, The stochastic finite element method: Past, present and future, *Comput Methods Appl Mech Eng*
537 198 (9-12) (2009) 1031–1051.
- 538 [9] M. Eiermann, O. G. Ernst, E. Ullmann, Computational aspects of the stochastic finite element method, in:
539 *Computing and Visualization in Science*, 2007.
- 540 [10] F. Tonon, A. Bernardini, A. Mammino, Reliability analysis of rock mass response by means of Random Set
541 Theory, *Reliab Eng Syst Saf* 70 (3) (2000) 263–282.
- 542 [11] I. M. Sobol, Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates,
543 *Math Comput Simul* 55 (1-3) (2001) 271–280.
- 544 [12] A. Saltelli, P. Annoni, I. Azzini, F. Campolongo, M. Ratto, S. Tarantola, Variance based sensitivity analysis
545 of model output. Design and estimator for the total sensitivity index, *Comput Phys Commun* 181 (2) (2010)
546 259–270.
- 547 [13] A. L. Eide, L. Holden, E. Reiso, S. I. Aanonsen, Automatic history matching by use of response surfaces and
548 experimental design, in: *ECMOR IV-4th European conference on the mathematics of oil recovery*, 1994.

- 549 [14] R. D. Hurrion, An example of simulation optimisation using a neural network metamodel: Finding the optimum
550 number of kanbans in a manufacturing system, *J Oper Res Soc* 48 (11) (1997) 1105–1112.
- 551 [15] G. I. Schuëller, H. J. Pradlwarter, P. S. Koutsourelakis, A critical appraisal of reliability estimation procedures
552 for high dimensions, *Probabilistic Eng Mech* 19 (4) (2004) 463–474.
- 553 [16] E. Zio, G. E. Apostolakis, Two methods for the structured assessment of model uncertainty by experts in
554 performance assessments of radioactive waste repositories, *Reliab Eng Syst Saf* 54 (2-3) (1996) 225–241.
- 555 [17] K. Sentz, S. Ferson, *Combination of Evidence in Dempster- Shafer Theory*, Vol. 4015, Citeseer, 2002.
- 556 [18] *Bounding uncertainty in civil engineering: Theoretical background and applications*, Vol. 539, Springer Science
557 & Business Media, 2012.
- 558 [19] H. F. Schweiger, G. M. Peschl, Reliability analysis in geotechnics with the random set finite element method,
559 *Comput Geotech* 32 (6) (2005) 422–435.
- 560 [20] C. F. Tsang, O. Stephansson, L. Jing, F. Kautsky, DECOVALEX Project: From 1992 to 2007, *Environ Geol*
561 57 (6) (2009) 1221–1237.
- 562 [21] 25 years of DECOVALEX - Scientific advances and lessons learned from an international research collaboration
563 in coupled subsurface processes, *Int J Rock Mech Min Sci* 122 (2019) 103995.
- 564 [22] S. Olivella, A. Gens, J. Carrera, E. E. Alonso, Numerical formulation for a simulator (CODE_BRIGHT) for
565 the coupled analysis of saline media, *Eng. Comput. (Swansea, Wales)* 13 (7) (1996) 87–112.
- 566 [23] J. Rutqvist, C. Tsang, TOUGH-FLAC: a numerical simulator for analysis of coupled thermal-hydrologic-
567 mechanical processes in fractured and porous geological media under multi-phase flow conditions, in: *Proceedings*
568 *of the TOUGH Symposium*, no. July, Lawrence Berkeley Natl. Lab. Berkeley, CA, 2003, pp. 1–9.
- 569 [24] O. Kolditz, S. Bauer, L. Bilke, N. Böttcher, J. O. Delfs, T. Fischer, U. J. Görke, T. Kalbacher, G. Kosakowski,
570 C. I. McDermott, C. H. Park, F. Radu, K. Rink, H. Shao, H. B. Shao, F. Sun, Y. Y. Sun, A. K. Singh, J. Taron,
571 M. Walther, W. Wang, N. Watanabe, Y. Wu, M. Xie, W. Xu, B. Zehner, OpenGeoSys: An open-source initiative
572 for numerical simulation of thermo-hydro-mechanical/chemical (THM/C) processes in porous media, *Environ*
573 *Earth Sci* 67 (2) (2012) 589–599.
- 574 [25] B. M. Adams, K. R. Dalbey, M. S. Eldred, D. M. Gay, L. P. Swiler, W. J. Bohnhoff, J. P. Eddy, P. D. Hough,
575 DAKOTA, A Multilevel Parallel Object Oriented Framework for Design Optimization, Parameter Estimation,
576 Uncertainty Quantification, and Sensitivity Analysis, Manual (December 2009) (2010) User Manual.
- 577 [26] E. Patelli, H. Pradlwarter, COSSAN-X : A General Purpose Code for Computational Stochastic Structural
578 Analysis COSSAN-X : A General Purpose Code for Computational Stochastic Structural Analysis Institute
579 of Engineering Mechanics University of Innsbruck, in: *IV European Conference on Computational Mechanics*,
580 Paris, France, EU, no. August 2018, 2010.
- 581 [27] X. Li, Entwicklung der Softwareplattform RESUS : Repository Simulation , Uncertainty propagation and Sensi-
582 tivity Analysis Dissertation, Clausthal University of Technology, Clausthal-Zellerfeld, Lower Saxony, Germany,
583 2015.
- 584 [28] D. Draper, A. Pereira, P. Prado, A. Saltelli, R. Cheal, S. Eguilior, B. Mendes, S. Tarantola, Scenario and para-
585 metric uncertainty in GESAMAC: A methodological study in nuclear waste disposal risk assessment, *Comput*
586 *Phys Commun* 117 (1) (1999) 142–155.
- 587 [29] A. Saltelli, S. Tarantola, On the relative importance of input factors in mathematical models: Safety assessment
588 for nuclear waste disposal, *J Am Stat Assoc* 97 (459) (2002) 702–709.
- 589 [30] S. J. Wang, K. C. Hsu, Stochastic Analysis of a Thermal Uncoupled Thermal-Hydraulic-Mechanical Model,
590 in: *Poromechanics 2017 - Proceedings of the 6th Biot Conference on Poromechanics*, American Society of Civil
591 Engineers (ASCE), 2017, pp. 787–794.
- 592 [31] L. Nguyen-Tuan, T. Lahmer, M. Datcheva, E. Stoimenova, T. Schanz, A novel parameter identification approach
593 for buffer elements involving complex coupled thermo-hydro-mechanical analyses, *Comput Geotech*.
- 594 [32] N. Watanabe, Finite element method for coupled thermo-hydro-mechanical processes in discretely fractured and
595 non-fractured porous media, Ph.D. thesis, Dresden University of Technology (2011).
- 596 [33] R. Lewis, B. Schrefler, *The Finite Element Method in the Static and Dynamic Deformation and Consolidation*
597 *of Porous Media*, 1st Edition, J. Wiley and Sons, 1998, an optional note.
- 598 [34] W. Ehlers, J. Bluhm, *Porous Media: Theory, Experiments and Numerical Applications*, 1st Edition, Springer,
599 2002.
- 600 [35] J. R. Booker, C. Savvidou, Consolidation around a point heat source., *Int J Numer Anal Methods Geomech*
601 9 (2) (1984) 173–184.
- 602 [36] A. A. Chaudhry, J. Buchwald, O. Kolditz, T. Nagel, Consolidation around a point heat source (correction and
603 verification), *Int J Numer Anal Methods Geomech* 43 (18) (2019) 2743–2751.
- 604 [37] H. R. Müller, B. Garitte, T. Vogt, S. Köhler, T. Sakaki, H. Weber, T. Spillmann, M. Hertrich, J. K. Becker,

- 605 N. Giroud, V. Cloet, N. Diomidis, T. Vietor, Implementation of the full-scale emplacement (FE) experiment at
606 the Mont Terri rock laboratory, in: *Swiss Journal of Geosciences*, Vol. 110, Springer, 2017, pp. 287–306.
- 607 [38] B. Li, E. W. Bhark, S. J. Gross, T. C. Billiter, K. Dehghani, Best practices of assisted history matching using
608 design of experiments, *SPE J* 24 (4) (2019) 1435–1451.
- 609 [39] C. Maschio, D. J. Schiozer, Bayesian history matching using artificial neural network and Markov Chain Monte
610 Carlo, *J Pet Sci Eng* 123 (2014) 62–71.
- 611 [40] C. Plúa, V. Manon, D. Seyedi, G. Armand, J. Rutqvist, J. Birkholzer, H. Xu, R. Guo, K. Tatcher, A. Bond,
612 W. Wang, T. Nagel, H. Shao, O. Kolditz, IBNL-2001265 (2020).
- 613 [41] M. Bossart, P. Thurry, Characteristics of the Opalinus Clay at Mont Terri, *Reports of the Swiss Geological
614 Survey no. 3* (3) (2008) 26.
- 615 [42] G. Armand, F. Bumbieler, N. Conil, R. de la Vaissière, J. M. Bosgiraud, M. N. Vu, Main outcomes from in
616 situ thermo-hydro-mechanical experiments programme to demonstrate feasibility of radioactive high-level waste
617 disposal in the Callovo-Oxfordian claystone, *J Rock Mech Geotech Eng* 9 (3) (2017) 415–427.
- 618 [43] R. L. Plackett, J. P. Burman, The Design of Optimum Multifactorial Experiments, *Biometrika* 33 (4) (1946)
619 305.
- 620 [44] GPy, GPy: A gaussian process framework in python (2014).
621 URL <http://github.com/SheffieldML/GPy>
- 622 [45] L. Bilke, B. Flemisch, T. Kalbacher, O. Kolditz, R. Helmig, T. Nagel, Development of open-source porous-media
623 simulators: principles and experiences., *Transp Porous Media*.
- 624 [46] S. Chatterjee, A. Hadi, *Regression Analysis by Example*, Wiley Series in Probability and Statistics, Wiley, 2006.
- 625 [47] O. Dubrule, Comparing splines and kriging, *Comput Geosci* 10 (2-3) (1984) 327–338.
- 626 [48] F. Xiong, Y. Xiong, W. Chen, S. Yang, Optimizing latin hypercube design for sequential sampling of computer
627 experiments, *Eng Optim* 41 (8) (2009) 793–810.