

Hydrogeological model selection among complex spatial priors

C. Brunetti¹, M. Bianchi², G. Pirot¹, N. Linde¹

¹Applied and Environmental Geophysics Group, Institute of Earth Sciences, University of Lausanne, 1015 Lausanne,

Switzerland

²British Geological Survey, Environmental Science Centre, Nottingham, UK

Key Points:

- Full Bayesian method for model selection among geologically-realistic conceptual models
- Thermodynamic integration and stepping-stone sampling provide consistent ranking of conceptual models
- Method applied to solute concentration data collected during a tracer test at the MADE site

Corresponding author: Carlotta Brunetti, Carlotta.Brunetti@unil.ch

Abstract

Hydrogeological field studies rely often on a single conceptual representation of the subsurface. This is problematic since the impact of a poorly chosen conceptual model on predictions might be significantly larger than the one caused by parameter uncertainty. Furthermore, conceptual models often need to incorporate geological concepts and patterns in order to provide meaningful uncertainty quantification and predictions. Consequently, several geologically-realistic conceptual models should ideally be considered and evaluated in terms of their relative merits. Here, we propose a full Bayesian methodology based on Markov chain Monte Carlo (MCMC) to enable model selection among 2D conceptual models that are sampled using training images and concepts from multiple-point statistics (MPS). More precisely, power posteriors for the different conceptual subsurface models are sampled using sequential geostatistical resampling and Graph Cuts. To demonstrate the methodology, we compare and rank five alternative conceptual geological models that have been proposed in the literature to describe aquifer heterogeneity at the MAcroDispersiOn Experiment (MADE) site in Mississippi, USA. We consider a small-scale tracer test (MADE-5) for which the spatial distribution of hydraulic conductivity impacts multilevel solute concentration data observed along a 2D transect. The thermodynamic integration and the stepping-stone sampling methods were used to compute the evidence and associated Bayes factors using the computed power posteriors. We find that both methods are compatible with MPS-based inversions and provide a consistent ranking of the competing conceptual models considered.

1 Introduction

The geological structure of the subsurface is a key controlling factor on groundwater flow and solute transport in aquifers [Maliva, 2016; Renard and Allard, 2013; Zheng and Gorelick, 2003] and, therefore, it needs to be properly represented and accounted for in modelling studies. The needs for quantitative and reliable subsurface modelling and management [Refsgaard and Henriksen, 2004; Scheidt et al., 2018] are driving hydrogeologists to consider conceptual models with increasing geological realism and complexity (e.g., see reviews by Linde et al. [2015a]; Hu and Chuginova [2008]). Traditionally, (hydro)geological subsurface heterogeneity has often been described in terms of mean values and covariances of the relevant physical properties (e.g., through the widely used multi-Gaussian models). However, such conceptualisations may be too simplistic in cer-

tain subsurface systems and, therefore, insufficient to accurately reproduce and predict flow and transport processes [Gómez-Hernández and Wen, 1998; Zinn and Harvey, 2003; Journel and Zhang, 2006; Kerrou *et al.*, 2008]. Multiple-point statistics (MPS) [Guardiano and Srivastava, 1993; Strebelle, 2002; Hu and Chugunova, 2008; Mariethoz and Caers, 2014] offers a means to effectively reproduce complex geological structures such as curvilinear features. By using a training image, MPS enables geostatistical simulations that honour point data and the higher-order spatial statistics that are captured in the training image. The training image is a conceptual representation summarising prior geological understanding about the system under study. It can be constructed from sketches drawn by hand, digitalised outcrops or generated by, for example, process-imitating, structure-imitating, or descriptive simulation methods [Koltermann and Gorelick, 1996; De Marsily *et al.*, 2005].

In many real world applications, generally because of the sparsity of direct observations, several alternative conceptualisations of subsurface heterogeneity (e.g., describing the spatial distribution of hydraulic conductivity) might be plausible and proposed by one or several experts. Unfortunately, uncertainty pertaining to the choice of the conceptual model is often ignored in modelling studies, even if it might be a dominant source of uncertainty [Bond *et al.*, 2007; Rojas *et al.*, 2008; Refsgaard *et al.*, 2012; Lark *et al.*, 2014; Scheidt *et al.*, 2018; Randle *et al.*, 2018]. Indeed, geostatistical model realisations generated from one training image might lead to a vastly different range of predictions than those generated from another training image, as shown, for example, by *Pirot et al.* [2015]. Conceptual uncertainty should, therefore, be integrated in modelling and inversion studies. Ideally, this should be achieved by using formal methods to test and rank alternative conceptual geological models based on available hydrogeological and geophysical data [Linde, 2014; Linde *et al.*, 2015a; Schöniger *et al.*, 2014; Dettmer *et al.*, 2010]. Bayesian model selection [Jeffreys, 1935, 1939; Kass and Raftery, 1995] offers a quantitative approach to perform such comparisons by computing the so called evidence (i.e., the denominator in Bayes' theorem) which allows to identify the conceptual model, in a chosen set, that is the most supported by the data. However, a complication arises when performing Bayesian model selection with complex spatial priors that are represented by training images. Most MPS-based inversions are non-parametric which implies that they rely on samples being drawn proportionally to the prior distribution, while it is generally not possible within a MPS framework to evaluate the prior probability of a given model proposal.

78 Hence, MPS-based inversions cannot build on many state-of-the-art concepts to enhance
79 the performance of the MCMC (e.g., *Laloy and Vrugt [2012]*) and associated approaches
80 for calculating the evidence [*Volpi et al., 2017; Brunetti et al., 2017*]. Similarly, it is not
81 possible within a MPS-framework to calculate approximate evidence estimates using the
82 Laplace-Metropolis method [*Lewis and Raftery, 1997*].

83 It is only recently that MPS-based inversions have been proposed (see review by
84 *Linde et al. [2015a]*). Markov chain Monte Carlo (MCMC) inversions with MPS (e.g., *Mariethoz et al. [2010a]; Hansen et al. [2012]*) generally rely on model proposals obtained
85 by sequential geostatistical resampling of the prior (Gibbs sampling) that are used within
86 the extended Metropolis algorithm to accept model proposals based on the likelihood ra-
87 tio [*Mosegaard and Tarantola, 1995*]. Sequential geostatistical resampling generates model
88 proposals of the spatially-distributed parameters of interest by conditional resimulations of
89 a random fraction of the current field proportionally to the prior as defined by the training
90 image. There exist several MPS methods to sample complex spatial priors with sequen-
91 tial Gibbs sampling. Examples include the versatile direct sampling method [*Mariethoz et al., 2010b*] or the recent Graph Cuts approach [*Zahner et al., 2016; Li et al., 2016*] that
92 enables speed-ups by one to two orders of magnitude. Since high-dimensional MCMC
93 inversions necessitate many evaluations of model proposals by forward modelling, it is
94 essential that the geostatistical model proposal process is fast compared to the forward
95 simulation time while ensuring model realisations of high quality that honour geological
96 patterns in the training image. Various advances have been made to enhance MPS-based
97 inversions both in a non-parametric MCMC framework (e.g., parallel tempering by *Laloy et al. [2016]*) and in a parametric framework using, for example, spatial generative adver-
98 sarial neural networks [*Laloy et al., 2018*]. Also, ensemble-based exploration schemes have
99 been explored [*Jäggli et al., 2017*].

103 State-of-the-art evidence estimators that are compatible with non-parametric spa-
104 tial priors include thermodynamic integration [*Gelman and Meng, 1998; Friel and Pettitt, 2008a*], stepping-stone [*Xie et al., 2011*] and nested sampling [*Skilling, 2004, 2006*]. The
105 thermodynamic integration method takes the name from its original application, which
106 was to compute the difference in a thermodynamic property (usually free energy) of a sys-
107 tem at two given states. Thermodynamic integration and the stepping-stone method sam-
108 ple from a sequence of so-called power posterior distributions that connect the prior to the
109 posterior distribution. The nested sampling method is based on a constrained local sam-
110

111 pling procedure in which the prior distribution is sampled under the constraint of a lower
112 bound on the log-likelihood function that increases with time. Thermodynamic integra-
113 tion and nested sampling transform the evidence, that is, a multi-dimensional integral over
114 the parameter space, into a one-dimensional integral over unit range in the log-likelihood
115 space. The stepping-stone sampling estimator approximates the evidence by importance
116 sampling using the power posteriors as importance distributions. To the best of our knowl-
117 edge, thermodynamic integration and stepping-stone sampling have never been used to
118 estimate the evidence of subsurface models built with MPS in the context of Bayesian
119 model selection, while this is the case for nested sampling [Elsheikh *et al.*, 2015]. Recent
120 studies in hydrology suggest that nested sampling is less accurate and stable than thermo-
121 dynamic integration [Liu *et al.*, 2016; Zeng *et al.*, 2018] and that it is strongly dependent
122 on the efficiency of the constrained local sampling procedure. Unfortunately, MPS-based
123 inversions cannot benefit from recent improvements in constrained local sampling ap-
124 proaches as they require parametric (analytical) forms of the prior [Schöniger *et al.*, 2014;
125 Liu *et al.*, 2016; Zeng *et al.*, 2018; Cao *et al.*, 2018]. Even if thermodynamic integration
126 and stepping-stone sampling are computationally expensive, they are easily parallelised
127 such that the computational time is equivalent to the time needed to run a single MCMC
128 chain. Moreover, these two methods are easy to implement and flexible in the sense that
129 any suitable MCMC method can, provided minimal changes, be used to explore the power
130 posterior distributions. The classical brute force Monte Carlo (MC) method [Hammersley
131 and Handscomb, 1964] can also be used to estimate the evidence when considering non-
132 parametric spatial priors. However, Brunetti *et al.* [2017] show that MC often requires a
133 prohibitive computational time to obtain reliable evidence estimates even for very simple
134 subsurface conceptualizations (e.g., layered models) when considering as few as seven un-
135 knowns. This limits its application to realistic high-dimensional MPS-based conceptual
136 models.

137 One way to circumvent the challenges of non-parametric priors in Bayesian model
138 selection is to reduce the model parameter space, for example, by cluster-based polynomial
139 chaos expansion [Bazargan and Christie, 2017] or by truncated discrete cosine transform
140 combined with summary metrics from training images [Lochbühler *et al.*, 2015]. Bayesian
141 inference and model selection is then applied on the reduced dimension space whose prior
142 distribution is parametric (e.g., multivariate Gaussian distribution). The main drawback

of such approaches is that truncation may smoothen sharp interfaces found in the training images.

In this study, we propose the first full Bayesian method that enables Bayesian model selection among geologically-realistic conceptual subsurface models. To do so, we combine sequential geostatistical resampling based on Graph Cuts, the extended Metropolis acceptance criterion and evidence estimation by power posteriors using either thermodynamic integration or stepping-stone sampling. The advantages and the drawbacks of this new methodology are assessed using a challenging application. In this study, we compare and rank five alternative conceptual geological models that have been proposed in the literature to characterise the spatial heterogeneity of the aquifer at the Macrodispersion Experiment (MADE) site in Mississippi, USA [Zheng *et al.*, 2011]. Among this set of five conceptual models of hydraulic conductivity spatial distribution, we aim to identify the one that is in the best agreement with multilevel concentration data acquired during a small-scale dipole tracer test (MADE-5) [Bianchi *et al.*, 2011a]. The case-study at the MADE site is used to demonstrate the ability of our Bayesian model selection method to deal with widely different conceptual hydrogeological models. We stress that the 2D modeling framework used herein limits our ability to generalize the findings to actual 3D field conditions. Extensions to 3D is methodologically straightforward, but computationally very challenging.

2 Theory

2.1 Bayesian inference and model selection

Bayesian inference approaches express the posterior pdf, $p(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$, of a set of unknown model parameters, $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_d\}$, given n measurements, $\tilde{\mathbf{Y}} = \{\tilde{y}_1, \dots, \tilde{y}_n\}$, via Bayes' theorem

$$p(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) = \frac{p(\boldsymbol{\theta})p(\tilde{\mathbf{Y}}|\boldsymbol{\theta})}{p(\tilde{\mathbf{Y}})}. \quad (1)$$

The prior pdf, $p(\boldsymbol{\theta})$, quantifies all the information that is available about the model parameters before considering the observed data. Typically, $p(\boldsymbol{\theta})$ is represented by multivariate analytical functions (e.g., Gaussian, uniform, exponential) describing marginal distributions of each parameter and their spatial correlation. With the advent of MPS methods, higher-order spatial statistics of $\boldsymbol{\theta}$ can be incorporated in inversions by means of training images. In this case, the description of prior knowledge is typically non-parametric and

173 sequential geostatistical resampling techniques are used to sample $p(\theta)$. The likelihood
 174 function, $p(\tilde{\mathbf{Y}}|\theta)$, summarises in a single scalar value the probability that the observed
 175 data has been generated by a proposed set of model parameters. We consider a Gaussian
 176 likelihood characterised by uncorrelated and normally distributed measurement errors with
 177 constant standard deviation, $\sigma_{\tilde{\mathbf{Y}}}$,

$$p(\tilde{\mathbf{Y}}|\theta) = \left(\sqrt{2\pi\sigma_{\tilde{\mathbf{Y}}}^2}\right)^{-n} \exp\left[-\frac{1}{2}\sum_{h=1}^n\left(\frac{\tilde{y}_h - \mathcal{F}_h(\theta)}{\sigma_{\tilde{\mathbf{Y}}}}\right)^2\right]. \quad (2)$$

178 As the residuals between the observed data, \tilde{y}_h , and the simulated forward responses,
 179 $\mathcal{F}_h(\theta)$, tends toward 0, the likelihood increases and, in particular, $p(\tilde{\mathbf{Y}}|\theta) \rightarrow \left(\sqrt{2\pi\sigma_{\tilde{\mathbf{Y}}}^2}\right)^{-n}$.
 180 The denominator in Bayes' theorem is the evidence (or marginal likelihood), $p(\tilde{\mathbf{Y}})$, and
 181 it is the cornerstone quantity in most Bayesian model selection problems. It should be
 182 noted, however, that the explicit computation of the evidence can be avoided by using re-
 183 versible jump (trans-dimensional) MCMC methods [Green, 1995]. The conceptual model
 184 with the highest evidence [Jeffreys, 1935, 1939] is the one that is the most supported by
 185 the data. A noteworthy feature of the evidence is that it implicitly accounts for the trade-
 186 off between goodness of fit and model complexity [Gull, 1988; Jeffreys, 1939; Jefferys and
 187 Berger, 1992; MacKay, 1992]. More precisely, the evidence quantifies how likely it is that
 188 a given conceptual model, $\eta \in \mathbb{N}$, with model parameters, θ , and prior distribution, $p(\theta|\eta)$,
 189 has generated the data $\tilde{\mathbf{Y}}$,

$$p(\tilde{\mathbf{Y}}|\eta) = \int p(\tilde{\mathbf{Y}}|\theta, \eta)p(\theta|\eta)d\theta. \quad (3)$$

190 The evidence is used to calculate Bayes factors [Kass and Raftery, 1995], that is, evidence
 191 ratios of one conceptual model with respect to an other. For instance, the Bayes factor of
 192 η_1 with respect to η_2 , or $B_{(\eta_1, \eta_2)}$, is defined as

$$B_{(\eta_1, \eta_2)} = \frac{p(\tilde{\mathbf{Y}}|\eta_1)}{p(\tilde{\mathbf{Y}}|\eta_2)}. \quad (4)$$

193 Conceptual models with large Bayes factors are preferred statistically and the conceptual
 194 model with the largest evidence is the one that best honours the data on average over its
 195 prior. However, the evidence computation is analytically intractable for most problems of
 196 interest and the multi-dimensional integral in Eq. 3 must be approximated by numerical
 197 means. In this work, the different conceptual models represent alternative spatial represen-
 198 tations of hydraulic conductivity in the subsurface.

2.2 Evidence estimation by power posteriors

Thermodynamic integration, also called path sampling [Gelman and Meng, 1998], and stepping-stone sampling [Xie et al., 2011] are two methods to estimate the evidence (Eq. 3) numerically. The key idea behind both methods is to sample from a sequence of so-called power posterior distributions, $p_\beta(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$, in order to create a path in the probability density space that connects the prior to the posterior distribution [Friel and Pettitt, 2008a]. The power posterior distribution is proportional to the prior pdf multiplied by the likelihood function raised to the power of $\beta \in [0, 1]$:

$$p_\beta(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto p(\boldsymbol{\theta})p(\tilde{\mathbf{Y}}|\boldsymbol{\theta})^\beta. \quad (5)$$

Decreasing β has the effect of flattening the likelihood function. For $\beta = 1$, the posterior distribution is sampled, $p_1(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto p(\boldsymbol{\theta})p(\tilde{\mathbf{Y}}|\boldsymbol{\theta})$; for $\beta = 0$, the prior distribution is sampled, $p_0(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto p(\boldsymbol{\theta})$. In thermodynamic integration and stepping-stone sampling, the priors are assumed to be proper and a sequence of β -values needs to be defined (see Section 2.2.3). For each β value, one (or more) MCMC runs are used to draw N samples from the corresponding power posterior distribution and the corresponding likelihood values are recorded. The Markov chains for the different β -values can be run independently in parallel or sequentially from $\beta = 0$ to $\beta = 1$ (serial MCMC) as described in Friel and Pettitt [2008a]. Thermodynamic integration and stepping-stone sampling have several attractive characteristics: (1) the total computing time is equivalent to a normal MCMC inversion provided that all MCMC runs are carried out in parallel, (2) they can be applied for any MCMC inversion method with only minimal intervention (it is only necessary to add the exponent β to the likelihood function) and (3) the only information needed is the series of likelihoods obtained from MCMC simulations with different β -values. Once the power posterior distributions have been sampled, the thermodynamic integration and stepping-stone sampling methods use the recorded likelihood values in two different ways to estimate the evidence (Sections 2.2.1-2.2.2).

2.2.1 Thermodynamic integration

Thermodynamic integration reduces the multi-dimensional integral of Eq. 3 into a one-dimensional integral of the expectation of the log-likelihood, $\log p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}, \eta)$, as:

$$\log p(\tilde{\mathbf{Y}}|\eta) = \int_0^1 E_{\boldsymbol{\theta}|\tilde{\mathbf{Y}},\beta} \left[\log p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}, \eta) \right] d\beta. \quad (6)$$

227 For the derivation of Eq. 6, we refer to *Friel and Pettitt [2008a]* and *Lartillot and Philippe*
 228 *[2006]*. The integral in Eq. 6 is estimated by a quadrature approximation over a discrete
 229 set of β -values, $0=\beta_1 < \dots < \beta_j < \dots < \beta_J=1$. To simplify the notation, we define the
 230 expectations of the log-likelihood functions as $\ell_j \equiv E_{\theta|\tilde{\mathbf{Y}},\beta_j} \left[\log p(\tilde{\mathbf{Y}}|\theta, \eta) \right]$ and their cor-
 231 responding variances as $\sigma_j^2 \equiv V_{\theta|\tilde{\mathbf{Y}},\beta_j} \left[\log p(\tilde{\mathbf{Y}}|\theta, \eta) \right]$. In this work, we use the corrected
 232 composite trapezoidal rule:

$$\log p(\tilde{\mathbf{Y}}|\eta) \approx \sum_{j=2}^J \frac{(\beta_j - \beta_{j-1})}{2} (\ell_j + \ell_{j-1}) - \sum_{j=2}^J \frac{(\beta_j - \beta_{j-1})^2}{12} (\sigma_j^2 - \sigma_{j-1}^2), \quad (7)$$

233 which provides more accurate estimates compared with the classical composite trapezoidal
 234 rule (first term in Eq. 7) as it also considers the second-order correction term (second
 235 term in Eq. 7). This corrected composite trapezoidal rule was originally employed by
 236 *Friel et al. [2014]* and later used by other authors including *Oates et al. [2016]* and *Grze-*
 237 *gorczyk et al. [2017]*.

238 The accuracy of the resulting evidence estimates depends on how the β -values are
 239 discretised, the number of β -values used, J , (details provided in Section 2.2.3), the num-
 240 ber, N , and the degree of correlation of the power posterior samples obtained by MCMC.
 241 The uncertainties associated with the evidence estimation by thermodynamic integration
 242 are often summarised by two error types: the sampling error, e_s , and the discretisation er-
 243 ror, e_d [*Lartillot and Philippe, 2006; Calderhead and Girolami, 2009*]. The sampling error
 244 is related to the standard errors of the MCMC posterior expectations of the log-likelihoods
 245 obtained for each β_j . To avoid underestimation of these errors, the autocorrelation in the
 246 MCMC samples should be accounted for in order to calculate the effective sample size,
 247 N_{eff} , (i.e., number of independent samples within each MCMC chain) as suggested by
 248 *Kass et al. [1998]*. The effective sample size is defined as:

$$N_{\text{eff},j} = \frac{N_j}{1 + 2 \sum_{z=1}^{\infty} \rho_j(z)}, \quad (8)$$

249 where $\rho_j(z)$ is the autocorrelation at lag z . Applying the rules for uncertainty propagation
 250 to the first leading term in Eq. 7 and assuming the errors of ℓ_j to be independent of those
 251 associated to ℓ_{j-1} , the sampling error is:

$$\sigma_s^2 = \sum_{j=2}^J \frac{(\beta_j - \beta_{j-1})^2}{4} \left(\frac{\sigma_j^2}{N_{\text{eff},j}} + \frac{\sigma_{j-1}^2}{N_{\text{eff},j-1}} \right). \quad (9)$$

252 Discretisation errors arise as the continuous integral of Eq. 6 is estimated using a finite
 253 number of evaluation points (Eq. 7). Following *Lartillot and Philippe [2006]*, *Baele et al.*
 254 *[2013]* and *Friel et al. [2014]*, we define e_d as the worst-case discretisation error that

255 arises from the approximation of Eq. 6 with a rectangular rule. Hence, e_d is half the dif-
 256 ference of the areas between the upper and lower step functions and it can be interpreted
 257 as the variance of the trapezoidal rule:

$$\sigma_d^2 = \sum_{j=2}^J \frac{(\beta_j - \beta_{j-1})^2}{4} (\ell_j - \ell_{j-1})^2. \quad (10)$$

258 As a consequence, the variance on the evidence estimates can be summarised as $\widehat{\text{Var}} \log p(\tilde{\mathbf{Y}}|\eta) =$
 259 $\sigma_d^2 + \sigma_s^2$.

260 2.2.2 Stepping-stone sampling

261 Stepping-stone sampling [Xie *et al.*, 2011] computes the evidence by combining
 262 power posteriors with importance sampling. The key underlying idea is to write the evi-
 263 dence as the ratio, r , of the normalising factors in Bayes' theorem for $\beta=1$ (posterior sam-
 264 pling) and $\beta=0$ (prior sampling):

$$r = \frac{p(\tilde{\mathbf{Y}}|\eta, \beta = 1)}{p(\tilde{\mathbf{Y}}|\eta, \beta = 0)}. \quad (11)$$

265 Since the prior integrates to one, the evidence is equivalent to r as $p(\tilde{\mathbf{Y}}|\eta, \beta = 0)$ equals 1.
 266 The ratio can be expressed as a product of J ratios, r_j :

$$r = \prod_{j=2}^J r_{j-1} = \prod_{j=2}^J \frac{p(\tilde{\mathbf{Y}}|\eta, \beta_j)}{p(\tilde{\mathbf{Y}}|\eta, \beta_{j-1})}. \quad (12)$$

267 Then, importance sampling is applied to the numerator and denominator of Eq. 12 using
 268 the power posterior $p_{\beta_{j-1}}(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$ as the importance distribution:

$$r_{j-1} = \frac{1}{N} \sum_{i=1}^N p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}_{j-1,i})^{\beta_j - \beta_{j-1}} \quad (13)$$

269 and, finally, the log-evidence is computed as:

$$\log p(\tilde{\mathbf{Y}}|\eta) = \sum_{j=2}^J \log r_{j-1} = \sum_{j=2}^J \log \left\{ \frac{1}{N} \sum_{i=1}^N \exp \left[(\beta_j - \beta_{j-1}) \cdot \log p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}_{j-1,i}) \right] \right\}. \quad (14)$$

270 In contrast to thermodynamic integration, the evidence estimated by stepping-stone sam-
 271 pling does not suffer from discretisation errors. The sampling error can be evaluated as:

$$\widehat{\text{Var}} \log p(\tilde{\mathbf{Y}}|\eta) = \sum_{j=2}^J \frac{1}{N_{\text{eff},j-1} \cdot N} \sum_{i=1}^N \left(\frac{p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}_{j-1,i})^{\beta_j - \beta_{j-1}}}{r_{j-1}} - 1 \right)^2. \quad (15)$$

272 The derivation of Eq. 14 and 15 appears in Xie *et al.* [2011]; Fan *et al.* [2011], and inter-
 273 ested readers are referred to this publication for further details. The only difference in our
 274 Eq. 15 is that we consider the effective sample size as defined in Eq. 8. Note that Eq. 13
 275 is only valid for the specific choice of $p_{\beta_{j-1}}(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$ as the importance distribution.

2.2.3 Discretisation scheme for β -values

For small increases of β close to 0, l_j increases dramatically and the corresponding power posteriors quickly turn from being similar to the prior to being similar to the posterior distribution (e.g., *Friel et al.* [2014]; *Oates et al.* [2016]; *Liu et al.* [2016]). As a consequence, the accuracy of the evidence estimates increases when placing most of the β -values close to 0 (e.g., *Friel and Pettitt* [2008b]; *Liu et al.* [2016]; *Grzegorzczuk et al.* [2017]). This is especially true for the thermodynamic integration method that estimates the evidence as the area below the curve of the expectation of the log-likelihood, l_j , as a function of β_j (Eq. 6). Starting from an initial set of sampling points, *Liu et al.* [2016] use an empirical method that places additional β -values based on a qualitative search for locations where l_j changes strongly in order to target additional β -values to use. However, this method is subjective and it increases the computing time when using parallel computations as the β -values are not defined at the outset. *Friel and Pettitt* [2008a] are the first to employ a discretisation scheme of β -values that follows a power law spacing as:

$$\beta_j = \left(\frac{j-1}{J-1} \right)^c \quad \text{with } j = 1, 2, \dots, J. \quad (16)$$

Calderhead and Girolami [2009] demonstrate that this scheme significantly improve the accuracy of the evidence estimates with respect to the uniform spacing used by *Lartillot and Philippe* [2006].

3 Method

3.1 General framework

It is common to sample the unnormalised posterior pdf of Eq. 1 with MCMC simulations. This is here achieved by combining the extended Metropolis acceptance criterion [*Mosegaard and Tarantola*, 1995] with a sequential geostatistical resampling technique (e.g., Graph Cuts) that provides conditional model proposals at each iteration featuring similar geological patterns as those found in the corresponding training image. For each proposed model, θ_{prop} , we calculate the forward response and compare it with the observed data and, according to the extended Metropolis algorithm, accept θ_{prop} with probability:

$$\alpha = \min \left\{ 1, \frac{p(\tilde{\mathbf{Y}}|\theta_{\text{prop}})}{p(\tilde{\mathbf{Y}}|\theta_{\text{cur}})} \right\}. \quad (17)$$

To sample the power posteriors, we simply modify the extended Metropolis acceptance criteria by raising the likelihoods in Eq. 17 with the corresponding β_k -values. We report

below the overall algorithm (Algorithm 1), in which we combine model proposals based on MPS with the extended Metropolis acceptance criteria followed by evidence estimation using power posteriors.

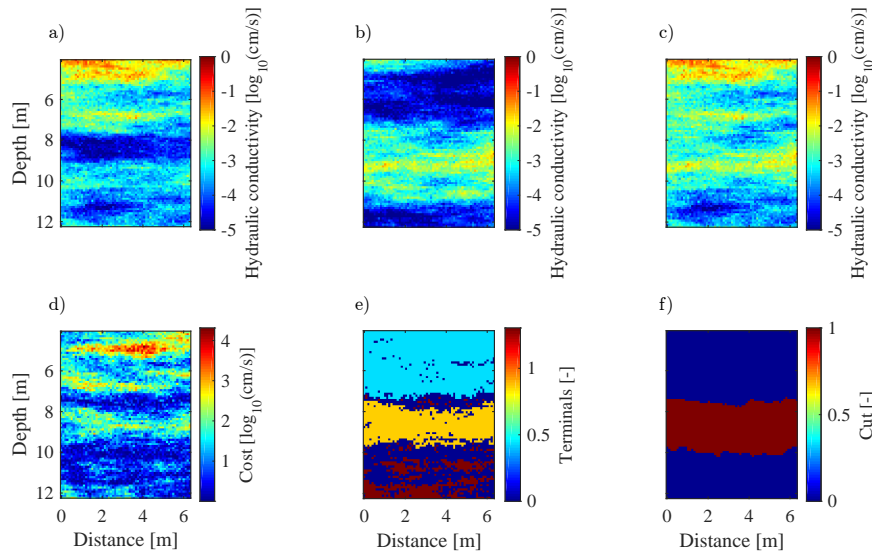
3.2 Graph Cuts model proposals

In this work, to sample spatially correlated parameters, we rely on model proposals based on the Graph Cuts algorithm introduced by *Zahner et al.* [2016] with some of the improvements proposed by *Pirot et al.* [2017a,b]. The main steps in the Graph Cuts algorithm are depicted in Figure 1. Basically, a section of the same size as the model domain, θ_{new} (Figure 1b), is randomly drawn from the training image and the absolute difference between θ_{new} and the current model realisation, θ_{cur} (Figure 1a), is computed and raised to the power of the cost power, δ_{cp} , [*Pirot et al.*, 2017b] to obtain the cost image, $\delta = |\theta_{\text{cur}} - \theta_{\text{new}}|^{\delta_{cp}}$ (Figure 1d). Two distinct regions of high cost, similar size and containing at least p pixels are randomly selected (Figure 1e). To choose these terminals, *Pirot et al.* [2017a] introduce the cutting threshold, $\delta_{th} \in [0, 100]$, defined as a percentile of $\max(\delta)$, which limits the possible terminals to those regions where $\delta > \delta_{th} \cdot \max(\delta)$. A patch is defined as the region enclosed by a minimum cost line separating the two terminals using the min-cut/max-flow algorithm by *Boykov and Kolmogorov* [2004] (Figure 1f) and the new model proposal, θ_{prop} (Figure 1c), is built by cutting the patch from θ_{new} and replacing the corresponding area in θ_{cur} .

We manually tune three algorithmic parameters to obtain model proposals that preserve the patterns found in the training image: the minimum number, p , of pixels in each of the two terminals, the cutting threshold, δ_{th} , and the cost power, δ_{cp} . We have set the cost power to 1 or 2 depending on the type of conceptual model considered. The main reason for using graph-cut proposals in this work is its computational speed relatively to other MPS algorithms (see comparisons by *Zahner et al.* [2016]). However, slower pixel-based geostatistical resimulation strategies that implement sequential Gibbs sampling, such as, those presented by *Mariethoz et al.* [2010b] or *Hansen et al.* [2012] could also be used.

3.3 Field site and available data

The MADE site is characterised by an unconsolidated shallow alluvial aquifer composed by a mixture of gravel, sand, and finer sediments. The high heterogeneity at the



324 **Figure 1.** Illustration of how model proposals are obtained using the Graph Cuts algorithm. (a) Current
 325 model realisation, θ_{cur} , (b) section drawn randomly from the training image, θ_{new} , and (c) the resulting
 326 model proposal, θ_{prop} . This model proposal is obtained as follows: (d) the cost image, δ , is defined as the
 327 absolute difference raised to the cost power, δ_{cP} , that is $\delta = |\theta_{\text{cur}} - \theta_{\text{new}}|^{\delta_{cP}}$, (e) two disconnected regions of
 328 high differences (light blue and orange areas) of similar size are randomly selected and (f) the cut of minimum
 329 cost that separates the two regions is calculated and the resulting dark red region is cut from (b) θ_{new} and
 330 pasted into (a) θ_{cur} to create (c) θ_{prop} .

342 MADE site got the attention of the hydrogeological community in the mid-1980s and nu-
 343 merous studies have been carried out since then (see *Zheng et al. [2011]* for a review).
 344 Previous interpretations of two large-scale tracer tests suggest that the structure is consis-
 345 tent with a network of highly permeable sediments embedded in a less permeable matrix
 346 [*Harvey and Gorelick, 2000; Feehley et al., 2000; Bianchi and Zheng, 2016*]. The case-
 347 study considered herein focuses on determining the most appropriate conceptual model
 348 of hydraulic conductivity in a reduced set given the multilevel solute concentration data
 349 collected during the MADE-5 tracer experiment [*Bianchi et al., 2011a*]. The test was per-
 350 formed in an array of four aligned boreholes with a maximum separation of 6 metres. The
 351 concentration data used in this work was collected in the two inner multi-level sampler
 352 (MLS) wells between the outer injection and abstraction wells, which were screened over
 353 the entire aquifer thickness. Before tracer injection, a steady-state dipole flow field was
 354 established by injecting clean water. Then, a known volume of bromide solution was in-
 355 jected along the entire vertical profile of the aquifer for 366 min followed by continuous

356 injection of clean water for 32 days. The flow rates at both the injection and extraction
357 wells were kept practically constant during all the steps of the test. Bromide concentra-
358 tions in the MLS wells were recorded at 19 different times and at seven depth levels (sam-
359 pling ports) in each of the two MLS wells resulting in 266 concentration measurements.
360 Full technical details about the experiment can be found in *Bianchi et al.* [2011a]. Given
361 the particular design of the borehole array, groundwater flow and bromide tracer trans-
362 port could be simulated only along the 2D transect intercepting the four wells (the forward
363 model used is described in Appendix A). This was necessary to reduce the computational
364 demands in this application of the proposed Bayesian model selection method. In prac-
365 tice, the 2D model assumes that the concentrations measured at the inner MLS wells are
366 mainly the result of transport along straight flow paths between the injection and the ab-
367 straction wells. To enable such 2D modeling, we performed a simple 3D-to-2D transfor-
368 mation of the data as described in Appendix A.

369 **3.3.1 Conceptual models at the MADE site and corresponding training images**

370 We consider five training images that may represent spatially distributed hydraulic
371 conductivity fields at the MADE site (Figure 2). The multi-Gaussian training image in
372 Figure 2a was created as a 2D unconditional realisation obtained with the Sequential Gaus-
373 sian SIMulation (SGSIM) algorithm of the Stanford Geostatistical Modeling Software
374 (SGeMS) [*Remy et al.*, 2009]. The corresponding variogram parameters (Table 1) were
375 calculated by *Bianchi et al.* [2011a] from the analysis of more than 1000 hydraulic con-
376 ductivity values estimated by means of borehole flowmeter tests [*Rehfeldt et al.*, 1992].
377 According to *Bianchi et al.* [2011a], the mean and variance in $\log_{10}(\text{cm/s})$ is set equal to
378 -2.37 and 1.95 , respectively.

389 The training images in Figure 2b-d were generated following *Linde et al.* [2015b].
390 The highly conductive and connected channels in an homogeneous matrix (Figure 2b)
391 is built from the original training image of *Strebelle* [2002] modified according to the
392 channel properties proposed by *Ronayne et al.* [2010] for the MADE site. The channel
393 hydraulic conductivity is equal to -0.54 in $\log_{10}(\text{cm/s})$, the channel thickness is 0.2 m and
394 the channel fraction is 3.25 %. The training image in Figure 2c is based on hydrogeo-
395 logical facies and their hydraulic conductivity values correspond to those of an outcrop
396 located near the MADE site [*Rehfeldt et al.*, 1992] and reported in Table 2.

399 The training image in Figure 2d is chosen solely on the knowledge that the aquifer
400 at the MADE site is constituted by alluvial deposits [Boggs *et al.*, 1992]. Linde *et al.* [2015b]
401 and Lochbühler *et al.* [2014] used the training image of Figure 2d as derived from a de-
402 tailed mapping study at the Herten site in Germany [Bayer *et al.*, 2011; Comunian *et al.*,
403 2011] featuring representative alluvial deposit structures and adapted it to the hydrogeo-
404 logical facies observed at the MADE site (Table 2).

405 The training image of Figure 2e is built based on five hydrogeological facies iden-
406 tified from lithological borehole data at the MADE site [Bianchi and Zheng, 2016] and
407 reported in Table 3. This training image is a stochastic unconditional realisation that was
408 generated following Bianchi and Zheng [2016].

411 Training images should be stationary and approach ergodicity [Caers and Zhang,
412 2004]. This implies that the type of patterns found should not change over the domain
413 covered by the training image (stationarity). Moreover, the size of the training image should
414 be sufficiently large (at least the double) compared to the largest pattern to enable ade-
415 quate simulations (ergodicity). Small training images lead to large ergodic fluctuations that
416 deteriorates pattern reproduction [Renard *et al.*, 2005]. Note that the smallest training im-
417 age considered herein (Figure 2b) is four times wider than the size of the model domain
418 in the horizontal direction.

419 In this work, we compare the five conceptual models of hydraulic conductivity that,
420 in the following, we refer to as (1) *multi-Gaussian* as built from the training image in Fig-
421 ure 2a; (2) *hybrid* that consists of the highly conductive channels of Figure 2b overlaid
422 on the multi-Gaussian background of Figure 2a; (3) *outcrop-based* built from the train-
423 ing image in Figure 2c; (4) *analog-based* built from the training image in Figure 2d; (5)
424 *lithofacies-based* built from the training image in Figure 2e. This selection of conceptual
425 models allows us to compare very different parameterisations of the spatial heterogene-
426 ity at the MADE site. Note that a full assessment of all conceptual models that has been
427 published for the MADE site is outside the scope of this study. Since computational lim-
428 itations prohibit full 3D simulations, we acknowledge that our findings in terms of the
429 suitability of different conceptual models at the MADE site should be treated with some
430 caution. Instead, the focus is on a new versatile methodology that enables comparison of
431 widely different conceptual models.

3.4 Evidence estimation in practice

We discretise the power coefficients β using the commonly used power law of Eq. 16 [Grzegorzczak *et al.*, 2017; Höhna *et al.*, 2017; Baele and Lemey, 2013; Xie *et al.*, 2011; Calderhead and Girolami, 2009; Friel and Pettitt, 2008a]. According to these studies, the parameter c should be set equal to 3 or 5 and J as large as possible with the common choice of $20 \leq J \leq 100$. In this study, we chose $c = 5$ and $J = 40$. For each β value, we run one MCMC chain of 10^5 iterations. These choices are dictated by computational constraints. The most challenging power posterior to sample is for $\beta=1$, for which we run 3 chains to better explore the posterior distribution. Consequently, we run 42 MCMC chains for each conceptual model. Given that the log-likelihoods obtained from the MCMC simulations are the basis for evidence estimations by power posteriors, we define the burn-in period (i.e., number of MCMC iterations required before reaching the target distribution) by considering the evolution of the log-likelihoods. To assess when the log-likelihood values start to oscillate around a constant value, we apply the Geweke method [Geweke, 1992] on the log-likelihoods of each chain. This diagnostic compares the mean computed on the last half of the considered chain length against the one derived from a smaller interval in the beginning of the chain (in our case, 20% of the chain length). At first, the Geweke's method is applied to the whole chain (no burn-in), and if its statistics is outside the 95% confidence interval of the standard normal distribution, we apply it again after discarding the first 1%, 2%, ..., 95% of the total chain length. The burn-in is determined in this way for $\beta=1$, as this is the most challenging case for which burn-in takes the longest time to achieve. The evidence estimates are computed using the thermodynamic integration method based on both the corrected trapezoidal rule (Eq. 7), as well as with the stepping-stone sampling method (Eq. 14). In order to correctly estimate the uncertainty of the evidence estimates, the effective sample size (Eq. 8) in each chain needs to be assessed. When evaluating Eq. 8, we truncate the sum in the denominator at the lag at which $\rho_j(z)$ is within 95% confidence interval of the normal distribution with standard deviation equal to the standard error of the sample autocorrelation. The evidence estimates are updated continuously after burn-in to visualise their evolution with the number of MCMC iterations. The uncertainty associated with the evidence estimates are summarised by standard errors, $SE = \sqrt{\widehat{\text{Var}} \log p(\tilde{\mathbf{Y}}|\eta)}$ with corresponding 95% confidence intervals. The variances $\widehat{\text{Var}} \log p(\tilde{\mathbf{Y}}|\eta)$ are computed using Eqs. 9-10 for the thermodynamic integration and using Eq. 15 for the stepping-stone sampling method.

4 Results for the MADE-5 case study

4.1 Bayesian inference

For each of the conceptual models considered, we first show prior MPS-realizations (i.e., $\beta = 0$) of hydraulic conductivity fields that are generated with the Graph Cuts method (Figure 3). Each set of prior realizations shows considerable spatial variability and is in broad agreement with the original training image (Figure 2). This is valid for both continuous (Figure 3b), categorical (Figures 3c-e) and hybrid conceptual models (Figure 3a).

The posterior distributions (i.e., $\beta = 1$) are obtained by assuming that the standard deviation of the measurement errors, $\sigma_{\bar{Y}}$ [mg/L], follows a log-uniform prior distribution in the range [1,10] mg/L (last column of Table 4). The lowest mean of the inferred $\sigma_{\bar{Y}}$ is obtained for the hybrid conceptual model (5.8 mg/L) suggesting that this model enables the best match with the data. The highest $\sigma_{\bar{Y}}$ is found for the outcrop-based model (9.4 mg/L). The acceptance rates are lower (second column in Table 4) than the ideal range between 15% and 40% proposed by *Gelman et al.* [1996], which suggests a slow convergence of the Markov chains. The burn-in time for each chain is obtained by the Geweke method (Table 4) as described in Section 3.4.

The different conceptual models provide quite different posterior distributions of the hydraulic conductivity field (Figure 4), even if certain commonalities are observed. For instance, all the posterior models have a high-conductive zone at a depth of 7 m that extends to a depth of 8 m on the right hand-side of the model domain. These features are visible in both the posterior mean and the maximum a-posteriori fields (first and second column of Figure 4). The analog- and outcrop-based conceptual models exhibit more variability in the inferred hydraulic conductivity values (Figures 4c and 4e) with respect to the others and the lithofacies-based conceptual model is characterised by the smallest posterior standard deviations (Figure 4d). The Gelman-Rubin statistic [*Gelman and Rubin*, 1992] is commonly used to assess if the MCMC chains has adequately sampled the posterior distribution, which is generally considered to be the case if this statistic is below 1.2. We see in the last column of Figure 4 that this is not the case for all pixel values, especially in the high-conductivity region, and that a larger number of iterations is required for a full convergence. However, we note that the evidence estimates are valid as long as the MCMC chains reach burn-in, while enhanced sampling decreases the estimation error.

509 In Figure 5, we show some of the simulated and observed breakthrough curves.
510 We have chosen the ones at a depth of 7 m in the monitoring wells MLS-1 (Figure 5a)
511 and MLS-2 (Figure 5b) because they correspond to a region of high conductivity (high
512 concentrations) and the ones at a depth of 11 m that correspond to low concentrations in
513 MLS-1 (Figure 5c) and MLS-2 (Figure 5d). Note that the range of measured concentration
514 values spans two orders of magnitude (Figure 5). In general, the outcrop-based concep-
515 tual model is the worst in reproducing the observed breakthrough curves while the hybrid
516 model is the best performing one; this is particularly clear in Figure 5d. Corresponding
517 plots at all measurement locations are found in the Supporting Information. The Pearson
518 correlation coefficients between the simulated posterior mean concentrations and the ob-
519 served ones are 0.96 for the hybrid model, 0.94 for the multi-Gaussian and analog-based
520 models, 0.91 for the lithofacies- and outcrop-based models.

525 **4.2 Bayesian model selection**

526 In this section, we present the estimated evidence values for each conceptual model
527 considered. Overall, the evidence values obtained using stepping-stone sampling and ther-
528 modynamic integration based on the corrected trapezoidal rule are in good agreement
529 with each other considering their 95% confidence intervals (Figure 6). Moreover, except
530 for some fluctuations at the early stage after burn-in, the evidence estimates evolve only
531 slowly as a function of the number of MCMC iterations after burn-in (Figure 6). We find
532 that stepping-stone sampling provides evidence values that are always lower than the ones
533 estimated with the thermodynamic integration. This behaviour is somewhat surprising as
534 the stepping-stone sampling technique is not based on a discretisation, while this is the
535 case for thermodynamic integration leading to an expected underestimation of the evi-
536 dence. The uncertainty associated with the stepping-stone evidence estimator decreases at
537 a sustained pace when increasing the number of MCMC iterations and it is lower than the
538 one associated with thermodynamic integration (Figure 6 and Table 5). Thermodynamic
539 integration is more affected by discretisation errors, an error source that is independent of
540 the number of MCMC iterations, than by sampling errors (Figure 8). For this reason, the
541 width of the confidence intervals obtained by thermodynamic integration does not reduce
542 significantly with increasing numbers of MCMC iterations (Figure 6).

550 Both evidence estimators lead to the same ranking of the conceptual models with
551 the hybrid conceptual model having the largest evidence and the outcrop-based conceptual

552 model having the lowest one (Table 5). The multi-Gaussian and the analog-based concep-
 553 tual models have very similar evidence estimates and they are the second-best performing
 554 conceptual models (Table 5).

558 For each conceptual model, the means of the log-likelihood functions, ℓ , increase
 559 with increasing β as we move from sampling the prior distribution ($\beta = 0$) to sampling
 560 the posterior distribution ($\beta = 1$) (Figure 7). From $\beta = 0$ to $\beta = 0.1$, the ℓ -estimates
 561 span three orders of magnitude. At very small values of β (i.e., $< 10^{-6}$), the outcrop-based
 562 conceptual model (green line in Figure 7) has mean log-likelihoods that are almost one
 563 order of magnitude higher than the other models. With increasing β , the outcrop-based
 564 model shows a much less steep increase of ℓ and at $\beta = 10^{-3}$, they start to be lower than
 565 the log-likelihood means of the other models. At higher power posteriors ($\beta > 0.1$), the
 566 ℓ -estimates for the hybrid conceptual model are the highest (red line in Figure 7), which
 567 explains why the highest evidence value is found for the hybrid conceptual model. We
 568 also note that the mean log-likelihood is not increasing continuously when β is close to
 569 one, which we attribute to random fluctuations of the MCMC chains (Figure 7).

572 The percentage ratio of independent MCMC samples after burn-in is never above
 573 10% and it decreases to values as low as 0.01% for $\beta = 1$ (Figure 8). This is a con-
 574 sequence of the slow mixing of the MCMC chains and it explains the increase of the
 575 sampling errors with increasing β for both thermodynamic integration (Figure 8c) and
 576 stepping-stone sampling (Figure 8d). The sampling errors of the stepping-stone sampling
 577 method are always at least two orders of magnitude higher than the ones related to the
 578 thermodynamic method, but this method is devoid of discretisation errors, which consti-
 579 tutes the dominant error type for thermodynamic integration. As mentioned before, using
 580 a power law to distribute β -values (Eq. 16) ensures that the discretisation errors for small
 581 β are relatively small (i.e., between 10^{-10} and 10^{-3} , Figure 8b). The pronounced fluctua-
 582 tions of the discretisation errors especially for $\beta > 0.1$ (Figure 8b) are related to the fact
 583 that the mean of the log-likelihoods does not increase monotonically for high β -values.

589 We now compute the Bayes factors for the best conceptual model (hybrid) with
 590 respect to each of the other competing conceptual models. In particular, we follow the
 591 guideline proposed by *Kass and Raftery* [1995] and we present twice the natural logarithm
 592 of the Bayes factors (Figures 9a-b). The Bayes factors of the hybrid conceptual model
 593 are on the order of 10^{15} and 10^{16} relative to the second best models (multi-Gaussian and

594 analog-based) and 10^{58} relative to the worst model (outcrop-based) for both thermody-
595 namic integration and stepping-stone sampling. Note that the measure of twice the natural
596 logarithms of the Bayes factors are all larger than 50 (Figures 9a-b). According to the in-
597 terpretation of *Kass and Raftery* [1995], we can safely claim that the hybrid model shows
598 very strong evidence of being superior to the other considered conceptual models. The
599 Bayes factors computed with the stepping-stone sampling method have smaller uncertain-
600 ties (Figure 9b) than the ones based on thermodynamic integration (Figure 9a). We note
601 that the relative rankings of the competing models obtained with the thermodynamic inte-
602 gration and the stepping-stone sampling methods are consistent and stable as long as the
603 MCMC chains has reached burn-in. Practically, this suggests that we can perform and ob-
604 tain reliable Bayesian model selection results at less computational cost and, again, that
605 formal convergence of the MCMC chains are not necessary.

611 5 Discussion

612 We have proposed a new methodology targeted at Bayesian model selection among
613 geologically-realistic conceptual models that are represented by training images. For MCMC
614 inversions, we use sequential geostatistical resampling through Graph Cuts that is two or-
615 ders of magnitude faster than the forward simulation time (i.e., 0.08 versus 8.35 sec). In
616 addition to being fast, the model realisations based on Graph Cuts are of high quality and
617 honour the geological patterns in the training images. This is true for the five different
618 types of conceptual models considered (Figures 3-4). Moreover, the Graph Cuts algorithm
619 can include point conditioning [*Li et al.*, 2016] even if this is not considered herein. In
620 our 2D analysis, we find that the hybrid conceptual model allows for the best fit of the ob-
621 served breakthrough curves (Figure 5). The inclusion of highly conductive channels in a
622 multi-Gaussian background enables enhanced simulations of the maximal concentrations
623 and it is in general agreement with the expected subsurface structure at the MADE site
624 (i.e., highly permeable network of sediments embedded in a less permeable matrix [*Har-*
625 *vey and Gorelick*, 2000; *Zheng and Gorelick*, 2003; *Liu et al.*, 2010; *Ronayne et al.*, 2010;
626 *Bianchi et al.*, 2011a,b]). We find that the outcrop model has not enough degrees of free-
627 dom to properly fit the solute concentration data (Figure 5). Furthermore, we expect that
628 an improved data fit would have been possible if we additionally would have inferred cer-
629 tain model parameter values (e.g., hydraulic conductivity for each facies and for the geo-
630 statistical parameters of the multi-Gaussian field).

631 In the light of the MADE-5 solute concentration data considered, the best fitting
 632 model (hybrid) is also the one that has the highest evidence, while the outcrop-based con-
 633 ceptual model has a Bayes factor of 10^{-58} with respect to the hybrid one, the lowest evi-
 634 dence and the lowest data fit (Table 4, Figure 6, Table 5). *Linde et al.* [2015b] rank differ-
 635 ent conceptual models (only the analog- and outcrop-based models are exactly the same as
 636 in the present work) of the region between the MLS-1 and MLS-2 wells using the maxi-
 637 mum likelihood estimate based on geophysical data (cross-hole ground-penetrating radar
 638 data). In agreement with our results, *Linde et al.* [2015b] find that the analog-based con-
 639 ceptual model explains the data much better than the outcrop-based conceptual model and
 640 that the latter is the worst performing one in the considered set.

641 Our results suggest that when comparing complex conceptual models represented by
 642 training images in data-rich environments, it may sometimes be possible to simply rank
 643 the performance of the competing conceptual models based on the inferred standard devi-
 644 ation of the measurement errors, $\sigma_{\bar{Y}}$ (Table 4), or the maximum likelihood estimate. Sim-
 645 ilar results for more traditional spatial priors were also found in other studies [*Schöniger*
 646 *et al.*, 2014; *Brunetti et al.*, 2017]. However, note that maximum likelihood-based model
 647 ranking will sometimes fail miserably as Bayesian model selection considers the trade-
 648 off between parsimony and goodness of fit. For instance, we expect that if we would have
 649 considered an uncorrelated hydraulic conductivity field, it would have produced the best
 650 fitting model but not the highest evidence. Moreover, it is also clear from these results
 651 that simply sampling the prior ($\beta = 0$) and then ranking the competing conceptual models
 652 based on the mean of the sampled likelihoods may provide misleading results. Indeed, the
 653 outcrop-based model has mean likelihoods of the prior model that are almost one order of
 654 magnitude higher than the ones of the other models (Figure 7) and, therefore, such a rank-
 655 ing would suggest that the outcrop-based conceptual model is the best one in describing
 656 the data while it is actually the worst one.

657 We find that stepping-stone sampling almost always provides slightly lower evi-
 658 dence estimates than thermodynamic integration (Table 5). This is in disagreement with
 659 the theory and with results by *Xie et al.* [2011] and *Friel et al.* [2014]. We attribute these
 660 unexpected results to the facts that (1) the MCMC chains for β close to 1 do not reach
 661 full convergence and the stepping-stone sampling is sensitive to poor convergence [*Friel*
 662 *et al.*, 2014] and (2) most of the contribution to the total evidence estimate comes from
 663 the terms of Eq. 7 computed for $\beta > 0.1$, a region where the mean log-likelihood does

664 not increase monotonically due to random fluctuations of the MCMC chains (Figure 7).
665 We also highlight that the comparison between the uncertainty estimates of the evidence
666 values provided by thermodynamic integration and stepping-stone sampling (Figure 6) is
667 not completely fair since the discretisation errors affecting thermodynamic integration are
668 based on a worst-case scenario that arises from the approximation of Eq. 6 with a rectan-
669 gular rule.

670 We stress again that our main intent is to present and demonstrate the proposed
671 methodology targeted at Bayesian model selection among geologically-realistic conceptual
672 models. Computational constraints made it infeasible to perform model selection in 3D.
673 Instead, given the particular design of the tracer experiment (i.e., array of four aligned
674 boreholes), we used a 2D flow and transport model and the data were corrected using
675 a 3D-to-2D transformation that account for differences in flowpaths for a homogeneous
676 subsurface (Appendix A). Since 3D heterogeneity is important at the MADE site, our 2D
677 model ranking should only be considered approximate.

678 Future work should better account for model errors caused by the 3D-to-2D flow
679 and transport approximation described in Appendix A. This would enhance the ability
680 to make more definite statements about aquifer heterogeneity at the MADE site. How to
681 properly account and represent model errors is a challenging task especially in problems
682 involving many data, high-dimensional parameter spaces and non-linear forward models
683 (e.g., *Linde et al. [2017]*). Another interesting topic that could be explored is to apply par-
684 allel tempering and use the resulting chains for computing the evidence with thermody-
685 namic integration or stepping-stone sampling [*Vlugt and Smit, 2001; Bailer-Jones, 2015;*
686 *Earl and Deem, 2005*]. Parallel tempering allows swapping between chains and, thereby,
687 improving sampling efficiency. This may contribute to more robust results, faster conver-
688 gence and, thereby, increase the number of effective samples (Figure 8a).

689 **6 Conclusions**

690 Inversions with geologically-realistic priors can be performed using training images
691 and model proposals that honour their multiple-point statistics. Unfortunately, such inver-
692 sions cannot rely on many state-of-the-art inversion methods and associated approaches for
693 calculating the evidence needed when performing Bayesian model selection. In this work,
694 we introduce a new full Bayesian methodology to enable Bayesian model selection among

695 complex geological priors. To demonstrate this methodology, we have evaluated its per-
696 formance in the context of determining, in a reduced set, the conceptual model that best
697 explains the concentration data for the case study considered (MADE-5). Our methodol-
698 ogy is applicable to both continuous and categorical conceptual models (e.g., a geologic
699 facies image) and it could be used at other sites, scales and for different data types. Ther-
700 modynamic integration and stepping-stone sampling methods are used for evidence com-
701 putation using a series of power posteriors obtained from MPS-based inversions. They
702 provide a consistent ranking of the competing conceptual models regardless of the number
703 of MCMC iterations after burn-in. This suggests that one can perform and obtain reliable
704 Bayesian model selection results with MCMC chains that have only achieved limited sam-
705 pling after burn-in. Both thermodynamic integration and stepping stone sampling are suit-
706 able evidence estimators. However, we recommend the stepping-stone sampling method
707 because it is not affected by discretisation errors and its uncertainty (sampling errors) is
708 significantly decreased with increasing numbers of MCMC iterations. This is not the case
709 for the thermodynamic integration because it is affected by discretisation errors that dom-
710 inate over the sampling errors. From the power posteriors derived from the test case, we
711 find that (1) ranking the conceptual models based on prior sampling only ($\beta = 0$) favours
712 the conceptual model with the lowest evidence and (2) model ranking based on the max-
713 imum posterior likelihood estimates ($\beta = 1$) provides, for this specific example, the same
714 results as the formal Bayesian model selection methods considered herein. For improved
715 sampling, we suggest that future work should investigate the use of parallel tempering re-
716 sults for evidence computations. Moreover, a full 3D analysis or a more formal treatment
717 of model errors due to the considered 3D-to-2D approximation would enhance the confi-
718 dence in statements about the suitability of alternative conceptual models at highly hetero-
719 geneous field sites.

720 **A: Forward model: from 3D to 2D**

721 The forward model used by *Bianchi et al.* [2011a] to simulate the bromide concen-
722 trations during the MADE-5 experiment is a 3D block-centred finite-difference model
723 based on MODFLOW (3D flow simulator) [*Harbaugh, 2005*] and MT3DMS (3D trans-
724 port simulator) [*Zheng, 2010*]. We initially consider a fine spatial discretisation of 0.1 m
725 in the area around the wells (Figure A.1a-b). However, running such a 3D model is com-
726 putationally prohibitive for evidence computations (i.e., 15 minutes of computing time to

727 get one forward response and we need 10^5 forward evaluations for each MCMC chain and
728 power posterior considered). To reduce the computing time, we perform a simple 3D to
729 2D correction of the data followed by 2D flow and transport simulations using the finite-
730 volume algorithm MaFloT [Künze and Lunati, 2012]. Moreover, we restrict the simulations
731 to the best fitting cross section (red segment in Figures A.1a-b) between the positions of
732 the injection, extraction and the two MLS wells, which results in an area of $6.3 \text{ m} \times 8.1$
733 m (Figure A.1c). For the transport equation, we set Dirichlet boundary conditions with the
734 normalised concentration to the given fluxes on the left side of the model domain (Fig-
735 ure A.1c) corresponding to the injection well location. For the pressure equation, we set
736 Dirichlet boundary conditions at the west and east sides (i.e., pressure difference), and
737 Neumann boundary conditions at the north and south sides of the model domain (Figure
738 A.1c).

746 Formal approaches to account for model errors in MCMC inversions exist (e.g., Cui
747 *et al.* [2011]), but they are outside the scope of the present contribution. In the following,
748 we introduce a simple error model that allows us to correct for the leading effects of the
749 3D to 2D transformation. These modelling errors stem primarily from the 2D linear ap-
750 proximation of the 3D radial distribution of the hydraulic heads, which results in a time
751 shift in the breakthrough curves at the MLS wells. To estimate the correction factors, we
752 consider a uniform hydraulic conductivity model with the geometric mean hydraulic con-
753 ductivity at the MADE site (i.e., $4.3 \cdot 10^{-5} \text{ m/s}$ [Rehfeldt *et al.*, 1992]). For this model, we
754 perform 3D and 2D simulations of the MADE-5 experiment with MODFLOW/MT3DMS
755 and MaFloT, respectively. As expected, the 3D simulated hydraulic heads between the
756 injection and extraction wells does not change linearly as for the 2D simulation (Figure
757 A.2). We tune the injection rate in the MODFLOW simulations to achieve simulated hy-
758 draulic heads that are as close as possible to the measured ones. We then perform MaFloT
759 simulations using the MODFLOW simulated hydraulic heads at the injection and extrac-
760 tion wells as boundary conditions and we compute correction factors at the MLS wells.
761 These multiplicative correction factors are those that maximise the correlation between
762 the concentrations simulated with MT3DMS and MaFloT. The mean correction factors
763 over the seven sampling ports in each of the two MLS wells are 1.09 and 1.92. Once the
764 correction factors have been applied, the earlier time shifts (Figures A.2b-c) are removed
765 (Figures A.2d-e). These correction factors are used in all subsequent simulations. Note
766 that no attempt is made to correct for tracer movement due to 3D heterogeneity; the cor-

767 rection is a simple geometrical correction to account for the transformation of a uniform
 768 3D to 2D flow field. We acknowledge that this is a crude approximation, but we deem it
 769 sufficient for the purposes of the present paper.

770 **Acknowledgments**

771 This work was supported by the Swiss National Science Foundation under grant num-
 772 ber 200021_155924. Niklas Linde thanks Arnaud Doucet for initially suggesting the use
 773 of thermodynamic integration. Marco Bianchi publishes with the permission of the Ex-
 774 ecutive Director of the British Geological Survey. The training images are available at
 775 <https://doi.org/10.5281/zenodo.2545587> and the concentration data of the MADE-5 tracer
 776 experiment will be soon available at <https://www.bgs.ac.uk/services/NGDC/>.

777 **References**

- 778 Baele, G., and P. Lemey (2013), Bayesian evolutionary model testing in the phyloge-
 779 nomics era: matching model complexity with computational efficiency, *Bioinformatics*,
 780 29(16), 1970–1979, doi:10.1093/bioinformatics/btt340.
- 781 Baele, G., P. Lemey, and S. Vansteelandt (2013), Make the most of your samples: Bayes
 782 factor estimators for high-dimensional models of sequence evolution, *BMC bioinformat-*
 783 *ics*, 14(1), 85, doi:10.1186/1471-2105-14-85.
- 784 Bailer-Jones, C. A. (2015), A general method for Bayesian time series modelling, *Tech.*
 785 *rep.*, Max Planck Institute for Astronomy, Heidelberg.
- 786 Bayer, P., P. Huggenberger, P. Renard, and A. Comunian (2011), Three-dimensional high
 787 resolution fluvio-glacial aquifer analog: Part 1: Field study, *Journal of Hydrology*,
 788 405(1-2), 1–9, doi:10.1016/j.jhydrol.2011.03.038.
- 789 Bazargan, H., and M. Christie (2017), Bayesian model selection for complex geological
 790 structures using polynomial chaos proxy, *Computational Geosciences*, 21(3), 533–551,
 791 doi:10.1007/s10596-017-9629-0.
- 792 Bianchi, M., and C. Zheng (2016), A lithofacies approach for modeling non-Fickian solute
 793 transport in a heterogeneous alluvial aquifer, *Water Resources Research*, 52(1), 552–565,
 794 doi:10.1002/2015WR018186.
- 795 Bianchi, M., C. Zheng, G. R. Tick, and S. M. Gorelick (2011a), Investigation of small-
 796 scale preferential flow with a forced-gradient tracer test, *Groundwater*, 49(4), 503–514,
 797 doi:10.1111/j.1745-6584.2010.00746.x.

- 798 Bianchi, M., C. Zheng, C. Wilson, G. R. Tick, G. Liu, and S. M. Gorelick (2011b), Spa-
799 tial connectivity in a highly heterogeneous aquifer: From cores to preferential flow
800 paths, *Water Resources Research*, 47(5), doi:10.1029/2009WR008966.
- 801 Boggs, J. M., S. C. Young, L. M. Beard, L. W. Gelhar, K. R. Rehfeldt, and E. E. Adams
802 (1992), Field study of dispersion in a heterogeneous aquifer: 1. Overview and site de-
803 scription, *Water Resources Research*, 28(12), 3281–3291, doi:10.1029/92WR01756.
- 804 Bond, C. E., A. D. Gibbs, Z. K. Shipton, and S. Jones (2007), What do you think this
805 is? "Conceptual uncertainty" in geoscience interpretation, *GSA today*, 17(11), 4, doi:
806 10.1130/GSAT01711A.1.
- 807 Boykov, Y., and V. Kolmogorov (2004), An experimental comparison of min-cut/max-flow
808 algorithms for energy minimization in vision, *IEEE transactions on pattern analysis and*
809 *machine intelligence*, 26(9), 1124–1137, doi:10.1109/TPAMI.2004.60.
- 810 Brunetti, C., N. Linde, and J. A. Vrugt (2017), Bayesian model selection in hydrogeo-
811 physics: Application to conceptual subsurface models of the South Oyster Bacte-
812 rial Transport Site, Virginia, USA, *Advances in Water Resources*, 102, 127–141, doi:
813 10.1016/j.advwatres.2017.02.006.
- 814 Caers, J., and T. Zhang (2004), Multiple-point geostatistics: a quantitative vehicle for in-
815 tegrating geologic analogs into multiple reservoir models, in *Integration of Outcrop*
816 *and Modern Analogs in Reservoir Modeling*, edited by G. M. Grammer, P. M. Har-
817 ris, and G. P. Eberli, chap. 18, American Association of Petroleum Geologists, doi:
818 10.1306/M80924C18.
- 819 Calderhead, B., and M. Girolami (2009), Estimating Bayes factors via thermodynamic
820 integration and population MCMC, *Computational Statistics & Data Analysis*, 53(12),
821 4028–4045, doi:10.1016/j.csda.2009.07.025.
- 822 Cao, T., X. Zeng, J. Wu, D. Wang, Y. Sun, X. Zhu, J. Lin, and Y. Long (2018), Integrat-
823 ing MT-DREAMzs and nested sampling algorithms to estimate marginal likelihood
824 and comparison with several other methods, *Journal of Hydrology*, 563, 750–765, doi:
825 10.1016/j.jhydrol.2018.06.055.
- 826 Comunian, A., P. Renard, J. Straubhaar, and P. Bayer (2011), Three-dimensional high res-
827 olution fluvio-glacial aquifer analog: Part 2: Geostatistical modeling, *Journal of hydrology*,
828 405(1-2), 10–23, doi:10.1016/j.jhydrol.2011.03.037.
- 829 Cui, T., C. Fox, and M. O'sullivan (2011), Bayesian calibration of a large-scale geothermal
830 reservoir model by a new adaptive delayed acceptance Metropolis Hastings algorithm,

- 831 *Water Resources Research*, 47(10), doi:10.1029/2010WR010352.
- 832 De Marsily, G., F. Delay, J. Gonçalves, P. Renard, V. Teles, and S. Violette (2005),
833 Dealing with spatial heterogeneity, *Hydrogeology Journal*, 13(1), 161–183, doi:
834 10.1007/s10040-004-0432-3.
- 835 Dettmer, J., S. E. Dosso, and J. C. Osler (2010), Bayesian evidence computation for model
836 selection in non-linear geoacoustic inference problems, *J Acoust Soc Am*, 128(6), 3406–
837 3415, doi:10.1121/1.3506345.
- 838 Earl, D. J., and M. W. Deem (2005), Parallel tempering: Theory, applications, and
839 new perspectives, *Physical Chemistry Chemical Physics*, 7(23), 3910–3916, doi:
840 10.1039/B509983H.
- 841 Elsheikh, A. H., V. Demyanov, R. Tavakoli, M. A. Christie, and M. F. Wheeler (2015),
842 Calibration of channelized subsurface flow models using nested sampling and soft prob-
843 abilities, *Advances in Water Resources*, 75, 14–30, doi:10.1016/j.advwatres.2014.10.006.
- 844 Fan, Y., R. Wu, M.-H. Chen, L. Kuo, and P. O. Lewis (2011), Choosing among partition
845 models in Bayesian phylogenetics, *Molecular biology and evolution*, 28(1), 523–532,
846 doi:10.1093/molbev/msq224.
- 847 Feehley, C. E., C. Zheng, and F. J. Molz (2000), A dual-domain mass transfer approach
848 for modeling solute transport in heterogeneous aquifers: Application to the Macrodis-
849 persion Experiment (MADE) site, *Water Resources Research*, 36(9), 2501–2515, doi:
850 10.1029/2000WR900148.
- 851 Friel, N., and A. N. Pettitt (2008a), Marginal likelihood estimation via power posteriors,
852 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(3), 589–
853 607, doi:10.1111/j.1467-9868.2007.00650.x.
- 854 Friel, N., and A. N. Pettitt (2008b), Marginal likelihood estimation via power posteriors,
855 *Journal of the Royal Statistical Society. Series B*, 70(3), 589–607, doi:10.1111/j.1467-
856 9868.2007.00650.x.
- 857 Friel, N., M. Hurn, and J. Wyse (2014), Improving power posterior estimation of statistical
858 evidence, *Statistics and Computing*, 24(5), 709–723, doi:10.1007/s11222-013-9397-1.
- 859 Gelman, A., and X.-L. Meng (1998), Simulating normalizing constants: From importance
860 sampling to bridge sampling to path sampling, *Statistical Science*, pp. 163–185, doi:
861 10.1214/ss/1028905934.
- 862 Gelman, A., and D. B. Rubin (1992), Inference from iterative simulation using multiple
863 sequences, *Statistical Science*, 7(4), 457–472, doi:10.1214/ss/1177011136.

- 864 Gelman, A., G. O. Roberts, and W. R. Gilks (1996), Efficient Metropolis jumping rules,
865 *Bayesian statistics*, 5, 599–608.
- 866 Geweke, J. (1992), Evaluating the accuracy of sampling-based approaches to the calcula-
867 tions of posterior moments, *Bayesian statistics*, 4, 641–649.
- 868 Gómez-Hernández, J. J., and X.-H. Wen (1998), To be or not to be multi-Gaussian? A
869 reflection on stochastic hydrogeology, *Advances in Water Resources*, 21(1), 47–61, doi:
870 10.1016/S0309-1708(96)00031-0.
- 871 Green, P. J. (1995), Reversible jump Markov chain Monte Carlo computation and Bayesian
872 model determination, *Biometrika*, 82(4), 711–732, doi:10.1093/biomet/82.4.711.
- 873 Grzegorzczak, M., A. Aderhold, and D. Husmeier (2017), Targeting Bayes factors with
874 direct-path non-equilibrium thermodynamic integration, *Computational Statistics*, 32(2),
875 717–761, doi:10.1007/s00180-017-0721-7.
- 876 Guardiano, F. B., and R. M. Srivastava (1993), Multivariate geostatistics: beyond bivariate
877 moments, in *Geostatistics Tróia '92*, edited by A. Soares, pp. 133–144, Springer, doi:
878 10.1007/978-94-011-1739-5_12.
- 879 Gull, S. F. (1988), Bayesian inductive inference and maximum entropy, in *Maximum-*
880 *entropy and Bayesian methods in Science and Engineering*, vol. 31-32, pp. 53–74,
881 Springer, doi:10.1007/978-94-009-3049-0_4.
- 882 Hammersley, J. M., and D. C. Handscomb (1964), *Monte Carlo methods*, vol. 1, VIII, 178
883 pp., Springer Netherlands, doi:10.1007/978-94-009-5819-7.
- 884 Hansen, T. M., K. S. Cordua, and K. Mosegaard (2012), Inverse problems with non-
885 trivial priors: Efficient solution through sequential Gibbs sampling, *Computational Geo-*
886 *sciences*, 16(3), 593–611, doi:10.1007/s10596-011-9271-1.
- 887 Harbaugh, A. W. (2005), *MODFLOW-2005, The US Geological Survey modular ground-*
888 *water model: the ground-water flow process*, US Department of the Interior, US Geolog-
889 ical Survey Reston.
- 890 Harvey, C., and S. M. Gorelick (2000), Rate-limited mass transfer or macrodispersion:
891 Which dominates plume evolution at the Macrodispersion Experiment (MADE) site?,
892 *Water Resources Research*, 36(3), 637–650, doi:10.1029/1999WR900247.
- 893 Höhna, S., M. L. Landis, and J. P. Huelsenbeck (2017), Parallel power posterior anal-
894 yses for fast computation of marginal likelihoods in phylogenetics, *bioRxiv*, doi:
895 doi.org/10.1101/104422.

- 896 Hu, L., and T. Chugunova (2008), Multiple-point geostatistics for modeling subsur-
897 face heterogeneity: A comprehensive review, *Water Resources Research*, 44(11), doi:
898 10.1029/2008WR006993.
- 899 Jäggli, C., J. Straubhaar, and P. Renard (2017), Posterior population expansion for
900 solving inverse problems, *Water Resources Research*, 53(4), 2902–2916, doi:
901 10.1002/2016WR019550.
- 902 Jefferys, W. H., and J. Berger (1992), Ockham’s Razor and Bayesian Analysis, *American*
903 *Scientist*, 80(1), 64–72.
- 904 Jeffreys, H. (1935), Some Tests of Significance, Treated by the Theory of Probability,
905 *Mathematical Proceedings of the Cambridge Philosophical Society*, 31(2), 203–222, doi:
906 10.1017/S030500410001330X.
- 907 Jeffreys, H. (1939), *Theory of Probability*, third ed., Oxford University Press.
- 908 Journal, A., and T. Zhang (2006), The necessity of a multiple-point prior model, *Mathe-*
909 *matical Geology*, 38(5), 591–610, doi:10.1007/s11004-006-9031-2.
- 910 Kass, R. E., and A. E. Raftery (1995), Bayes factors, *Journal of the American Statistical*
911 *Association*, 90(430), 773–795, doi:10.1080/01621459.1995.10476572.
- 912 Kass, R. E., B. P. Carlin, A. Gelman, and R. M. Neal (1998), Markov chain Monte Carlo
913 in practice: a roundtable discussion, *The American Statistician*, 52(2), 93–100, doi:
914 10.1080/00031305.1998.10480547.
- 915 Kerrou, J., P. Renard, H.-J. H. Franssen, and I. Lunati (2008), Issues in characterizing het-
916 erogeneity and connectivity in non-multigaussian media, *Advances in Water Resources*,
917 31(1), 147–159, doi:10.1016/j.advwatres.2007.07.002.
- 918 Koltermann, C. E., and S. M. Gorelick (1996), Heterogeneity in sedimentary deposits:
919 A review of structure-imitating, process-imitating, and descriptive approaches, *Water*
920 *Resources Research*, 32(9), 2617–2658, doi:doi.org/10.1029/96WR00025.
- 921 Künze, R., and I. Lunati (2012), An adaptive multiscale method for density-
922 driven instabilities, *Journal of Computational Physics*, 231(17), 5557–5570, doi:
923 10.1016/j.jcp.2012.02.025.
- 924 Laloy, E., and J. A. Vrugt (2012), High-dimensional posterior exploration of hydrologic
925 models using multiple-try DREAM_{ZS} and high-performance computing, *Water Re-*
926 *sources Research*, 48(1), 1–18, doi:10.1029/2011WR010608.
- 927 Laloy, E., N. Linde, D. Jacques, and G. Mariethoz (2016), Merging parallel tempering
928 with sequential geostatistical resampling for improved posterior exploration of high-

- 929 dimensional subsurface categorical fields, *Advances in Water Resources*, 90, 57–69, doi:
930 10.1016/j.advwatres.2016.02.008.
- 931 Laloy, E., R. Héroult, D. Jacques, and N. Linde (2018), Training-image based geostatistical
932 inversion using a spatial generative adversarial neural network, *Water Resources
933 Research*, 54(1), 381–406, doi:10.1002/2017WR022148.
- 934 Lark, R., S. Thorpe, H. Kessler, and S. Mathers (2014), Interpretative modelling of a geological
935 cross section from boreholes: sources of uncertainty and their quantification,
936 *Solid Earth*, 5(2), 1189–1203, doi:10.5194/se-5-1189-2014.
- 937 Lartillot, N., and H. Philippe (2006), Computing Bayes factors using thermodynamic inte-
938 gration, *Systematic Biology*, 55(2), 195–207, doi:10.1080/10635150500433722.
- 939 Lewis, S. M., and A. E. Raftery (1997), Estimating Bayes factors via posterior simulation
940 with the Laplace-Metropolis estimator, *J Am Stat Assoc*, 92(438), 648–655, doi:
941 10.1080/01621459.1997.10474016.
- 942 Li, X., G. Mariethoz, D. Lu, and N. Linde (2016), Patch-based iterative conditional geo-
943 statistical simulation using graph cuts, *Water Resources Research*, 52(8), 6297–6320,
944 doi:10.1002/2015WR018378.
- 945 Linde, N. (2014), Falsification and corroboration of conceptual hydrological models
946 using geophysical data, *Wiley Interdisciplinary Reviews: Water*, 1(2), 151–171, doi:
947 10.1002/wat2.1011.
- 948 Linde, N., P. Renard, T. Mukerji, and J. Caers (2015a), Geological realism in hydrogeo-
949 logical and geophysical inverse modeling: A review, *Advances of Water Resources*, 86,
950 86–101, doi:10.1016/j.advwatres.2015.09.019.
- 951 Linde, N., T. Lochbühler, M. Dogan, and R. L. Van Dam (2015b), Tomogram-based
952 comparison of geostatistical models: Application to the Macrodispersion Experiment
953 (MADE) site, *Journal of Hydrology*, 531, 543–556, doi:10.1016/j.jhydrol.2015.10.073.
- 954 Linde, N., D. Ginsbourger, J. Irving, F. Nobile, and A. Doucet (2017), On Uncertainty
955 Quantification in Hydrogeology and Hydrogeophysics, *Advances in Water Resources*,
956 110, 166–181, doi:10.1016/j.advwatres.2017.10.014.
- 957 Liu, G., C. Zheng, G. R. Tick, J. J. Butler, and S. M. Gorelick (2010), Relative impor-
958 tance of dispersion and rate-limited mass transfer in highly heterogeneous porous me-
959 dia: Analysis of a new tracer test at the Macrodispersion Experiment (MADE) site, *Wa-
960 ter Resources Research*, 46(3), doi:10.1029/2009WR008430.

- 961 Liu, P., A. S. Elshall, M. Ye, P. Beerli, X. Zeng, D. Lu, and Y. Tao (2016), Evaluating
962 marginal likelihood with thermodynamic integration method and comparison with
963 several other numerical methods, *Water Resources Research*, 52(2), 734–758, doi:
964 10.1002/2014WR016718.
- 965 Lochbühler, T., G. Pirot, J. Straubhaar, and N. Linde (2014), Conditioning of multiple-
966 point statistics facies simulations to tomographic images, *Mathematical Geosciences*,
967 46(5), 625–645, doi:10.1007/s11004-013-9484-z.
- 968 Lochbühler, T., J. A. Vrugt, M. Sadegh, and N. Linde (2015), Summary statistics from
969 training images as prior information in probabilistic inversion, *Geophys J Int*, 201(1),
970 157–171, doi:10.1093/gji/ggv008.
- 971 MacKay, D. J. (1992), Bayesian interpolation, *Neural Computation*, 4(3), 415–447, doi:
972 10.1162/neco.1992.4.3.415.
- 973 Maliva, R. G. (2016), *Aquifer Characterization Techniques*, Springer, doi:10.1007/978-3-
974 319-32137-0.
- 975 Mariethoz, G., and J. Caers (2014), *Multiple-point Geostatistics: Stochastic Modeling with*
976 *Training Images*, John Wiley & Sons, doi:10.1002/9781118662953.
- 977 Mariethoz, G., P. Renard, and J. Caers (2010a), Bayesian inverse problem and opti-
978 mization with iterative spatial resampling, *Water Resources Research*, 46(11), doi:
979 doi.org/10.1029/2010WR009274.
- 980 Mariethoz, G., P. Renard, and J. Straubhaar (2010b), The direct sampling method to per-
981 form multiple-point geostatistical simulations, *Water Resources Research*, 46(11), doi:
982 10.1029/2008WR007621.
- 983 Mosegaard, K., and A. Tarantola (1995), Monte Carlo sampling of solutions to inverse
984 problems, *Journal of Geophysical Research: Solid Earth*, 100(B7), 12,431–12,447, doi:
985 10.1029/94JB03097.
- 986 Oates, C. J., T. Papamarkou, and M. Girolami (2016), The controlled thermodynamic inte-
987 gral for Bayesian model evidence evaluation, *Journal of the American Statistical Associ-*
988 *ation*, 111(514), 634–645, doi:10.1080/01621459.2015.1021006.
- 989 Pirot, G., P. Renard, E. Huber, J. Straubhaar, and P. Huguenberger (2015), Influence of
990 conceptual model uncertainty on contaminant transport forecasting in braided river
991 aquifers, *Journal of Hydrology*, 531, 124–141, doi:10.1016/j.jhydrol.2015.07.036.
- 992 Pirot, G., N. Linde, G. Mariethoz, and J. H. Bradford (2017a), Probabilistic inversion with
993 graph cuts: Application to the Boise Hydrogeophysical Research Site, *Water Resources*

- 994 *Research*, 53(2), 1231–1250, doi:doi.org/10.1002/2016WR019347.
- 995 Pirot, G., M. Cardiff, G. Mariethoz, J. Bradford, and N. Linde (2017b), Towards 3D Prob-
996 abilistic Inversion with Graphcuts, in *23rd European Meeting of Environmental and En-
997 gineering Geophysics*.
- 998 Randle, C. H., C. E. Bond, R. M. Lark, and A. A. Monaghan (2018), Can uncertainty in
999 geological cross-section interpretations be quantified and predicted?, *Geosphere*, doi:
1000 10.1130/GES01510.1.
- 1001 Refsgaard, J. C., and H. J. Henriksen (2004), Modelling guidelines—terminology
1002 and guiding principles, *Advances in Water Resources*, 27(1), 71–82, doi:
1003 10.1016/j.advwatres.2003.08.006.
- 1004 Refsgaard, J. C., S. Christensen, T. O. Sonnenborg, D. Seifert, A. L. Højberg, and
1005 L. Trolborg (2012), Review of strategies for handling geological uncertainty in ground-
1006 water flow and transport modeling, *Advances in Water Resources*, 36, 36–50, doi:
1007 10.1016/j.advwatres.2011.04.006.
- 1008 Rehfeldt, K. R., J. M. Boggs, and L. W. Gelhar (1992), Field study of dispersion in a het-
1009 erogeneous aquifer: 3. Geostatistical analysis of hydraulic conductivity, *Water Resources
1010 Research*, 28(12), 3309–3324, doi:10.1029/92WR01758.
- 1011 Remy, N., A. Boucher, and J. Wu (2009), *Applied geostatistics with SGeMS: a user's guide*,
1012 Cambridge University Press.
- 1013 Renard, P., and D. Allard (2013), Connectivity metrics for subsurface flow and transport,
1014 *Advances in Water Resources*, 51, 168–196, doi:10.1016/j.advwatres.2011.12.001.
- 1015 Renard, P., H. Demougeot-Renard, and R. Froidevaux (2005), *Geostatistics for environ-
1016 mental applications*, Springer, doi:10.1007/s11004-018-9733-2.
- 1017 Rojas, R., L. Feyen, and A. Dassargues (2008), Conceptual model uncertainty in ground-
1018 water modeling: Combining generalized likelihood uncertainty estimation and Bayesian
1019 model averaging, *Water Resources Research*, 44(12), doi:10.1029/2008WR006908.
- 1020 Ronayne, M. J., S. M. Gorelick, and C. Zheng (2010), Geological modeling of submeter
1021 scale heterogeneity and its influence on tracer transport in a fluvial aquifer, *Water Re-
1022 sources Research*, 46(10), doi:10.1029/2010WR009348.
- 1023 Scheidt, C., L. Li, and J. Caers (2018), *Quantifying Uncertainty in Subsurface Systems*, vol.
1024 236, John Wiley & Sons, doi:10.1002/9781119325888.
- 1025 Schöniger, A., T. Wöhling, L. Samaniego, and W. Nowak (2014), Model selection on solid
1026 ground: Rigorous comparison of nine ways to evaluate Bayesian model evidence, *Water*

- 1027 *Resources Research*, 50(12), 9484–9513, doi:10.1002/2014WR016062.
- 1028 Skilling, J. (2004), Nested sampling, in *AIP Conference Proceedings*, vol. 735, pp. 395–
1029 405, AIP, doi:10.1063/1.1835238.
- 1030 Skilling, J. (2006), Nested sampling for general Bayesian computation, *Bayesian analysis*,
1031 1(4), 833–859, doi:10.1214/06-BA127.
- 1032 Strebelle, S. (2002), Conditional simulation of complex geological structures
1033 using multiple-point statistics, *Mathematical Geology*, 34(1), 1–21, doi:
1034 10.1023/A:1014009426274.
- 1035 Vlugt, T. J., and B. Smit (2001), On the efficient sampling of pathways in the transition
1036 path ensemble, *PhysChemComm*, 4(2), 11–17, doi:10.1039/B009865P.
- 1037 Volpi, E., G. Schoups, G. Firmani, and J. Vrugt (2017), Sworn testimony of the model
1038 evidence: Gaussian Mixture Importance (GAME) sampling, *Water Resources Research*,
1039 doi:10.1002/2016WR020167.
- 1040 Xie, W., P. O. Lewis, Y. Fan, L. Kuo, and M.-H. Chen (2011), Improving marginal like-
1041 lihood estimation for Bayesian phylogenetic model selection, *Systematic Biology*, 60(2),
1042 150–160, doi:10.1093/sysbio/syq085.
- 1043 Zahner, T., T. Lochbühler, G. Mariethoz, and N. Linde (2016), Image synthesis with graph
1044 cuts: a fast model proposal mechanism in probabilistic inversion, *Geophysical Journal*
1045 *International*, 204(2), 1179–1190, doi:10.1093/gji/ggv517.
- 1046 Zeng, X., M. Ye, J. Wu, D. Wang, and X. Zhu (2018), Improved nested sampling and
1047 surrogate-enabled comparison with other marginal likelihood estimators, *Water Re-*
1048 *sources Research*, 54(2), 797–826, doi:10.1002/2017WR020782.
- 1049 Zheng, C. (2010), MT3DMS v5. 3 supplemental user's guide, *Department of Geological*
1050 *Sciences, University of Alabama, Tuscaloosa, Alabama*.
- 1051 Zheng, C., and S. M. Gorelick (2003), Analysis of solute transport in flow fields influ-
1052 enced by preferential flowpaths at the decimeter scale, *Groundwater*, 41(2), 142–155,
1053 doi:10.1111/j.1745-6584.2003.tb02578.x.
- 1054 Zheng, C., M. Bianchi, and S. M. Gorelick (2011), Lessons learned from 25 years
1055 of research at the MADE site, *Groundwater*, 49(5), 649–662, doi:10.1111/j.1745-
1056 6584.2010.00753.x.
- 1057 Zinn, B., and C. F. Harvey (2003), When good statistical models of aquifer heterogeneity
1058 go bad: A comparison of flow, dispersion, and mass transfer in connected and multi-
1059 variate Gaussian hydraulic conductivity fields, *Water Resources Research*, 39(3), doi:

1060

10.1029/2001WR001146.

Algorithm 1: MCMC inversion workflow based on MPS and the extended Metropolis

algorithm to enable evidence estimation using power posteriors.

Input: T , maximum number of MCMC iterations; J , number of power coefficients β

distributed according to Eq. 16; a training image

Output: Λ_j , matrices containing power posteriors and log-likelihoods; $\log p(\tilde{\mathbf{Y}}|\eta)$,

evidence

Set $t = 1$;

Draw θ_1 from the training image;

Solve the forward problem;

Compute likelihood (e.g., Eq. 2);

for $j = 1, \dots, J$ **do**

for $t = 2, \dots, T$ **do**

 Set $\theta_{\text{cur}} = \theta_{t-1}$;

 Draw θ_{prop} based on MPS (e.g., using Graph Cuts proposals);

 Solve the forward problem;

 Compute likelihood (e.g., Eq. 2);

 Accept θ_{prop} with probability, $\alpha = \min \left\{ 1, \frac{p(\tilde{\mathbf{Y}}|\theta_{\text{prop}})^{\beta_j}}{p(\tilde{\mathbf{Y}}|\theta_{\text{cur}})^{\beta_j}} \right\}$;

if θ_{prop} *accepted* **then**

 Set $\theta_t = \theta_{\text{prop}}$;

else

 Set $\theta_t = \theta_{\text{cur}}$;

end

 Store θ_t and the corresponding log-likelihood in matrix Λ_j ;

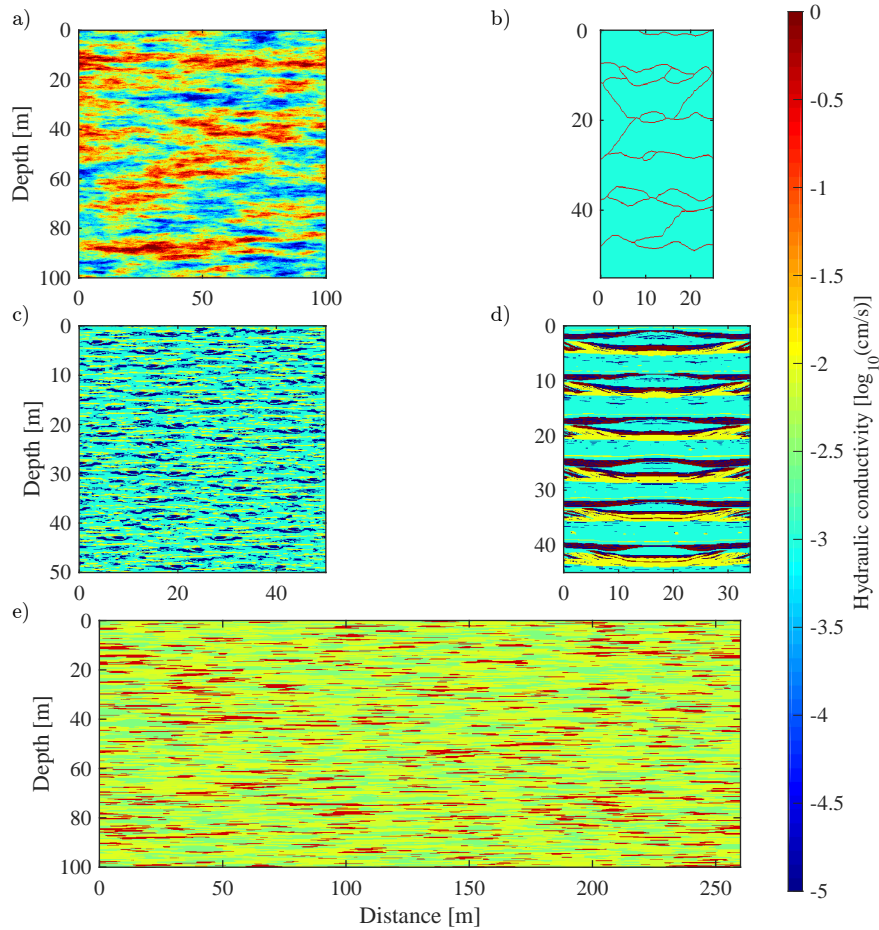
 Set $t=t+1$;

end

end

Compute $\log p(\tilde{\mathbf{Y}}|\eta)$ (Eqs. 7 and 14) and corresponding variances (Eqs. 9-10 and 15)

using the information stored in Λ_j after the removal of the burn-in period.



379 **Figure 2.** Training images used in the MPS-based inversion to represent spatial hydraulic conductivity of
 380 the MADE site: (a) multi-Gaussian field [Bianchi *et al.*, 2011a], (b) highly conductive channels in an homoge-
 381 neous matrix [Strebelle, 2002; Ronayne *et al.*, 2010; Linde *et al.*, 2015b], (c) model based on a mapping study
 382 of a MADE outcrop [Rehfeldt *et al.*, 1992; Linde *et al.*, 2015b], (d) model based on a mapping study at the
 383 Herten site in Germany [Bayer *et al.*, 2011; Comunian *et al.*, 2011; Linde *et al.*, 2015b] featuring representa-
 384 tive alluvial deposit structures and (e) model based on lithological borehole data collected at the MADE site
 385 [Bianchi and Zheng, 2016].

386 **Table 1.** Geostatistical parameters of the multi-Gaussian training image (Figure 2a) proposed by *Bianchi*
 387 *et al.* [2011a] for the MADE site. The actual variogram model was a linear combination of a spherical and an
 388 exponential model.

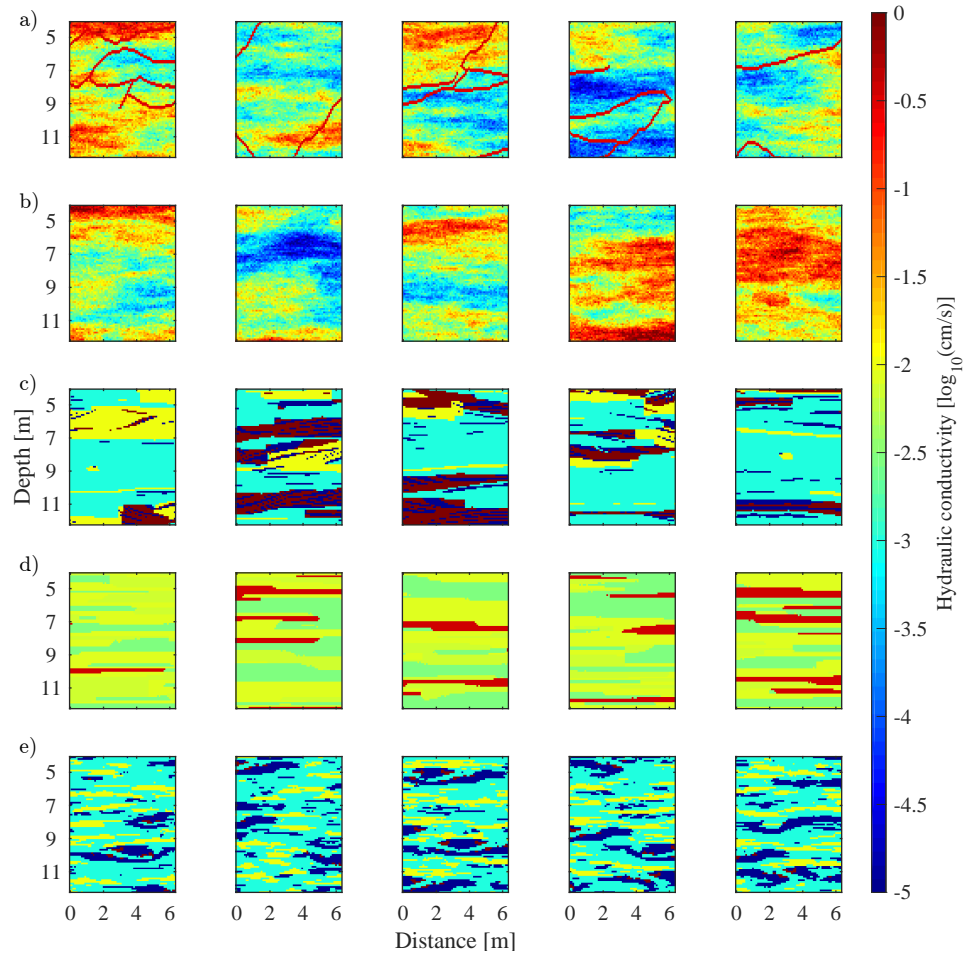
Variogram parameters	Variogram model	
	Spherical	Exponential
Maximum range [m]	76	21
Minimum range [m]	4.6	5
Nugget	0.2	-
Sill	1.75	3.0

397 **Table 2.** Hydrogeological facies and their hydraulic conductivity values [*Rehfeldt et al.*, 1992] observed at
 398 the MADE site outcrop and used for the training images in Figure 2c-d.

Facies	$\log_{10} K$ [cm/s]
Open framework gravel	$-6.83 \cdot 10^{-4}$
Sand	-2.00
Undifferentiated sandy gravel	-3.00
Sandy, clayey gravel	-5.00

409 **Table 3.** Hydrogeological facies and their hydraulic conductivity values based on lithological data from the
 410 MADE site [*Bianchi and Zheng*, 2016] and used for the training image in Figure 2e.

Facies	$\log_{10} K$ [cm/s]
Highly conductive gravel	-0.45
Sand and gravel	-2.05
Gravel with sand	-2.11
Well-sorted sand	-2.18
Sand gravel and fines	-2.53



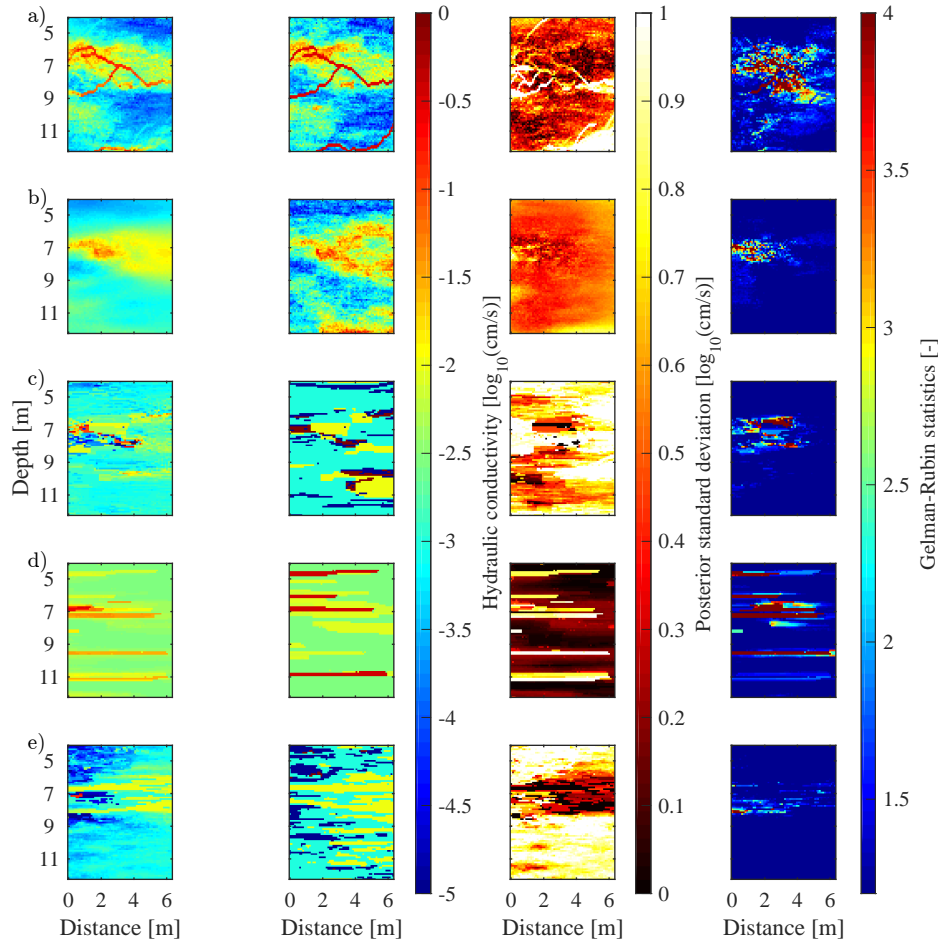
472 **Figure 3.** Five prior realisations of hydraulic conductivity fields generated from the training images of Fig-
 473 ure 2 with the Graph Cuts algorithm for the (a) hybrid, (b) multi-Gaussian, (c) analog-based, (d) lithofacies-
 474 based and (e) outcrop-based conceptual model of the MADE site.

484 **Table 4.** Summary of MCMC results using the MADE-5 tracer data for three MCMC chains of 10^5 steps
 485 for each conceptual model with $\beta = 1$. First column, conceptual model considered; second column, average
 486 acceptance rate (AR); third to fifth column, burn-in percentage based on the Geweke method for each of the
 487 three chains (when no value is displayed, the chain failed to reach burn-in); last two columns, means and
 488 standard deviations of the standard deviation of the measurement errors inferred with MCMC.

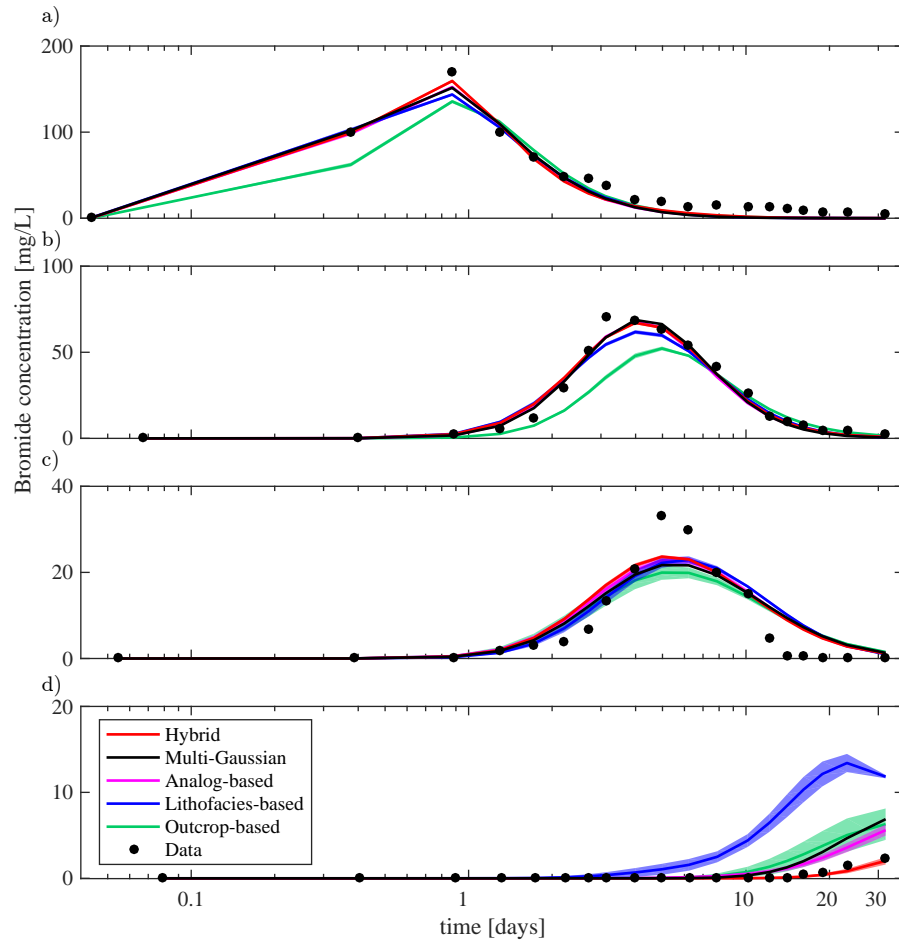
Conceptual model	AR [%]	Burn-in [%]			$\sigma_{\tilde{Y}}$ [mg/L]	
		Chain 1	Chain 2	Chain 3	Mean	Std
Hybrid	0.6	-	58	87	5.81	0.27
Multi-Gaussian	8.0	48	45	62	7.14	0.33
Analog	4.1	-	64	84	7.22	0.34
Lithofacies	1.2	55	38	74	8.92	0.60
Outcrop	5.5	76	97	-	9.36	0.35

555 **Table 5.** Estimates of the natural logarithm of the evidence, $\log p(\tilde{Y}|\eta)$, with corresponding standard errors,
 556 SE, for each conceptual model (first column) based on the stepping-stone sampling method (second and third
 557 column) and on the thermodynamic integration method with the corrected trapezoidal rule (last two columns).

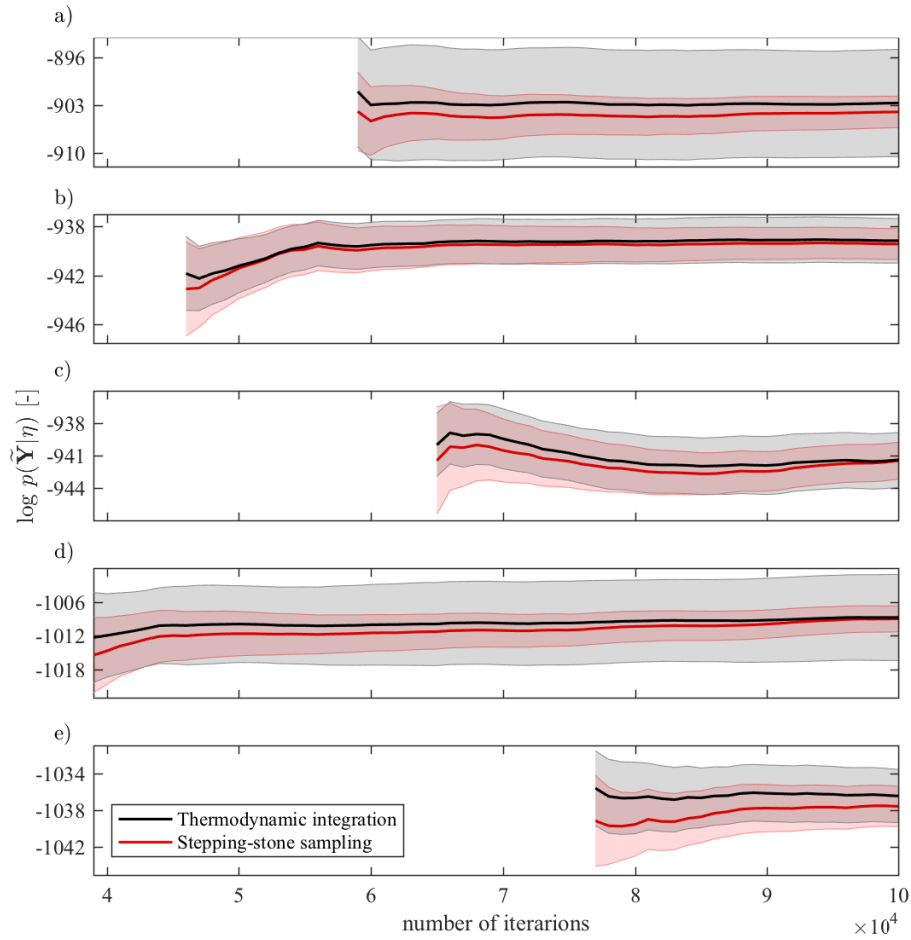
Conceptual model	Stepping-stone sampling		Thermodynamic integration	
	$\log p(\tilde{Y} \eta)$ [-]	SE [-]	$\log p(\tilde{Y} \eta)$ [-]	SE [-]
Hybrid	-903.99	1.17	-902.68	4.02
Multi-Gaussian	-939.43	0.64	-939.15	0.93
Analog	-941.48	0.87	-941.40	1.30
Lithofacies	-1009.01	1.18	-1008.76	3.90
Outcrop	-1037.58	1.11	-1036.45	1.47



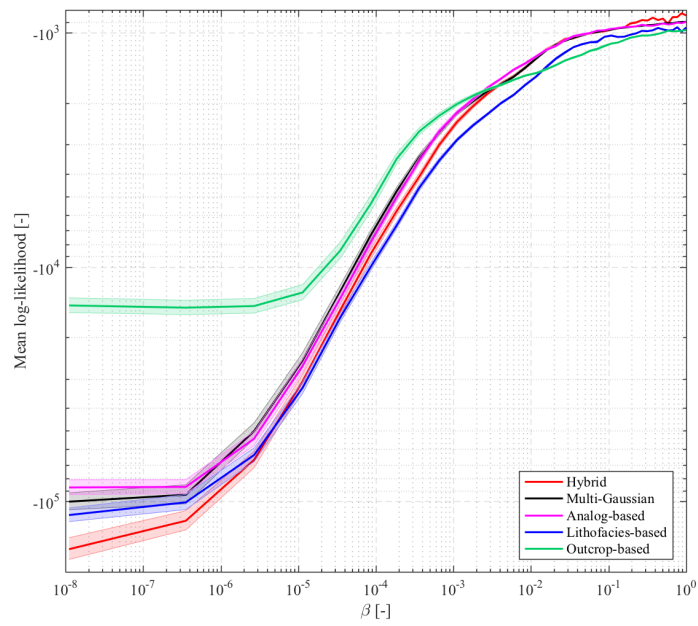
504 **Figure 4.** Mean (first column), maximum a-posteriori (second column), and standard deviation (third
 505 column) of the posterior hydraulic conductivity realisations for the (a) hybrid, (b) multi-Gaussian, (c) analog-
 506 based, (d) lithofacies-based and (e) outcrop-based conceptual model at the MADE site. In the last column, the
 507 Gelman-Rubin statistic for each pixel is reported. Dark-blue regions represent values equal or less than 1.2
 508 and indicate that convergence has been reached for those pixels.



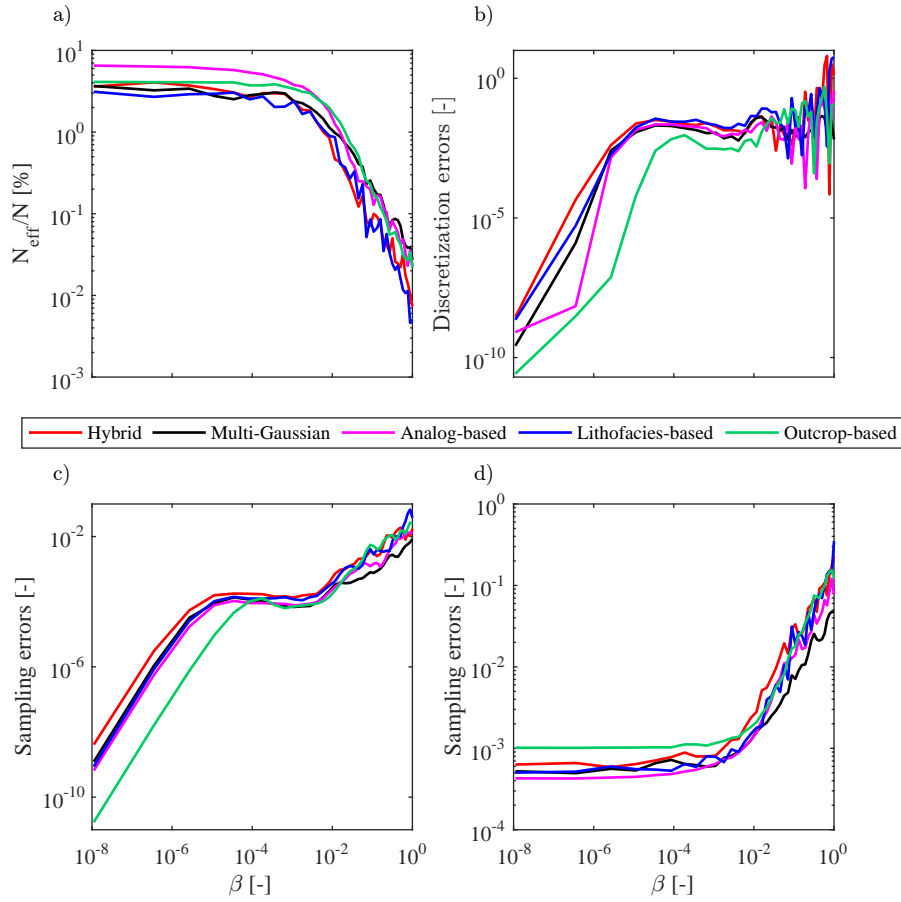
521 **Figure 5.** Simulated (solid lines) and measured (black dots) bromide breakthrough curves from the MADE-
 522 5 experiment in the two monitoring wells MLS-1 and MLS-2 at a depth of 7 m (a-b) and 11 m (c-d), respec-
 523 tively. The simulated breakthrough curves are summarised by the mean of the posterior realisations (solid
 524 lines) and their 95% confidence intervals (shaded areas).



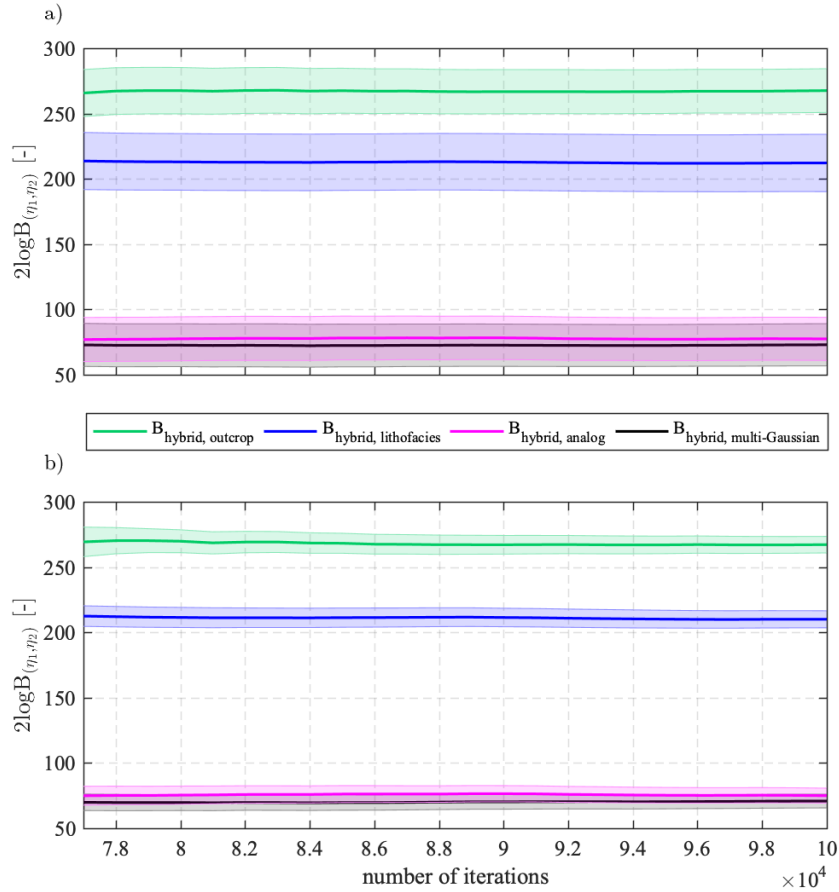
543 **Figure 6.** Natural logarithm of the evidence estimates, $\log p(\tilde{\mathbf{Y}}|\eta)$, as a function of the number of MCMC
 544 iterations. Evidences are computed with the stepping-stone sampling method (red line) and the thermo-
 545 dynamic integration method based on the corrected trapezoidal rule (black line) for the (a) hybrid, (b)
 546 multi-Gaussian, (c) analog-based, (d) lithofacies-based and (e) outcrop-based model at the MADE site.
 547 The evidence computation starts after a different number of MCMC iterations because each of the conceptual
 548 models has a specific burn-in period. The shaded areas represent the 95% confidence interval of the evidence
 549 estimates (pink for stepping-stone sampling and grey for thermodynamic integration).



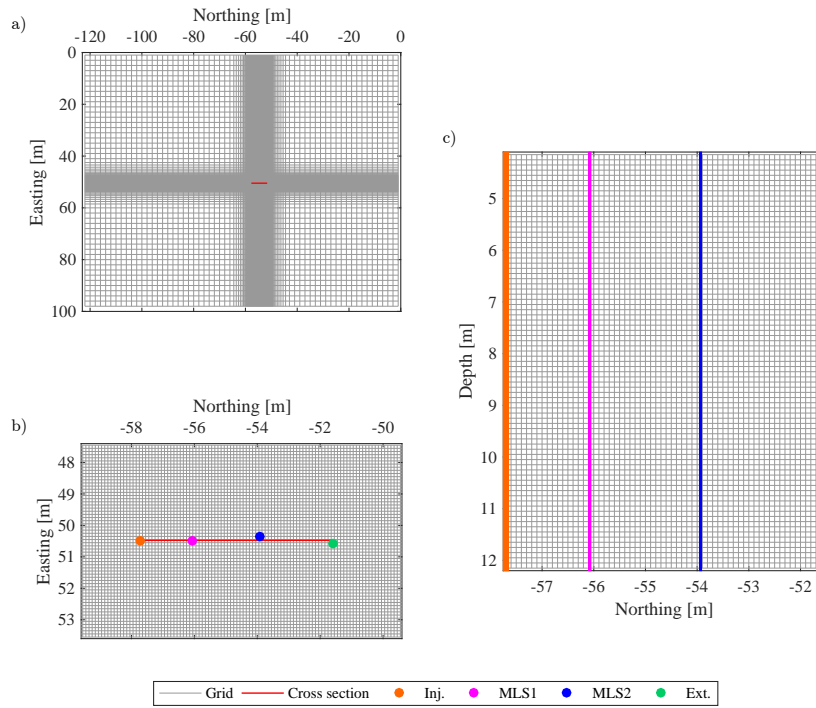
570 **Figure 7.** Mean (line) of the natural logarithm of the likelihood functions, ℓ , computed for each β value and
 571 the 95% confidence interval of the ℓ -estimates (shaded areas). Note that the x - and y -axes are in \log_{10} scale.



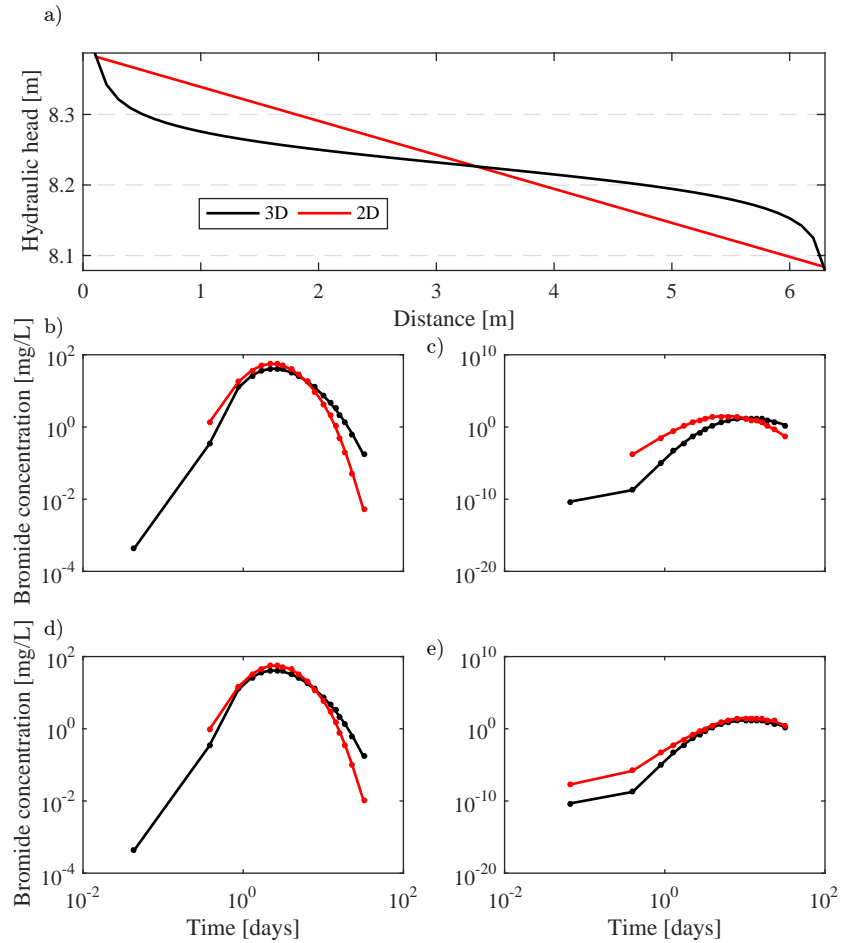
584 **Figure 8.** (a) Percentage ratio between the effective and the total number of MCMC samples, (b) discreti-
 585 cation errors in the thermodynamic integration method (square root of Eq. 10), (c) sampling errors in the
 586 thermodynamic integration method (square root of Eq. 9) and (d) sampling errors in the stepping-stone sam-
 587 pling method (square root of Eq. 15) as a function of β -values. Note that all the x - and y -axes are in \log_{10}
 588 scale.



606 **Figure 9.** Twice the natural logarithm of the Bayes factors of the "best model" (hybrid) with respect to the
 607 outcrop-based (green line), lithofacies-based (blue line), analog-based (magenta line) and multi-Gaussian
 608 (black line) conceptual model at the MADE site. Results are shown for (a) the thermodynamic integration
 609 method based on the corrected trapezoidal rule and for the (b) stepping-stone sampling method. The shaded
 610 areas represent the 95% confidence interval of the $2\log B_{\eta_1, \eta_2}$ measures.



739 **Figure A.1.** (a) Aerial view of the 3D grid used for simulations with MODFLOW/MT3DMS; (b) zoom
 740 in the tracer test area, in which the grid size was refined to 0.1 m; (c) cross section used for simulations with
 741 MaFloT. The width of the lines in (c) is representative of the diameter of the four wells.



742 **Figure A.2.** (a) Hydraulic head profiles between the injection and extraction wells arising from 2D and 3D
 743 flow simulations in a uniform hydraulic conductivity field. Simulated breakthrough curves at 7 m depth in (b)
 744 MLS-1 and (c) MLS-2 without corrections. The shifts in the 2D simulations are removed when (d-e) applying
 745 the correction factors.