

# Merging parallel tempering with sequential geostatistical resampling for improved posterior exploration of high-dimensional subsurface categorical fields



Eric Laloy<sup>a,\*</sup>, Niklas Linde<sup>b</sup>, Diederik Jacques<sup>a</sup>, Grégoire Mariethoz<sup>c</sup>

<sup>a</sup> Institute for Environment, Health and Safety, Belgian Nuclear Research Centre, Belgium

<sup>b</sup> Applied and Environmental Geophysics Group, Institute of Earth Sciences, University of Lausanne

<sup>c</sup> Institute of Earth Surface Dynamics, University of Lausanne, Lausanne, Switzerland

## ARTICLE INFO

### Article history:

Received 5 November 2015

Revised 16 February 2016

Accepted 16 February 2016

Available online 24 February 2016

### Keywords:

Parallel tempering

Sequential geostatistical resampling

Training image

MCMC

Multiple-point statistics

## ABSTRACT

The sequential geostatistical resampling (SGR) algorithm is a Markov chain Monte Carlo (MCMC) scheme for sampling from possibly non-Gaussian, complex spatially-distributed prior models such as geologic facies or categorical fields. In this work, we highlight the limits of standard SGR for posterior inference of high-dimensional categorical fields with realistically complex likelihood landscapes and benchmark a parallel tempering implementation (PT-SGR). Our proposed PT-SGR approach is demonstrated using synthetic (error corrupted) data from steady-state flow and transport experiments in categorical 7575- and 10,000-dimensional 2D conductivity fields. In both case studies, every SGR trial gets trapped in a local optima while PT-SGR maintains a higher diversity in the sampled model states. The advantage of PT-SGR is most apparent in an inverse transport problem where the posterior distribution is made bimodal by construction. PT-SGR then converges towards the appropriate data misfit much faster than SGR and partly recovers the two modes. In contrast, for the same computational resources SGR does not fit the data to the appropriate error level and hardly produces a locally optimal solution that looks visually similar to one of the two reference modes. Although PT-SGR clearly surpasses SGR in performance, our results also indicate that using a small number (16–24) of temperatures (and thus parallel cores) may not permit complete sampling of the posterior distribution by PT-SGR within a reasonable computational time (less than 1–2 weeks).

© 2016 Elsevier Ltd. All rights reserved.

## 1. Introduction

A general Markov chain Monte Carlo (MCMC) simulation strategy based on sequential geostatistical resampling of spatially-distributed prior models has recently been introduced in the geosciences to infer the posterior distribution of subsurface property fields. The approach creates candidate fields by conditioning a geostatistical field realization from a possibly complex prior model to a randomly chosen fraction of the current state (and hence model/field) of the Markov chain. Both parametric (e.g., multi-Gaussian) and non-parametric prior models can be considered. The multi-Gaussian prior basically consists of a variogram model that

encodes the 2-point statistics to be honored. As of non-Gaussian structures, they can be generated using a multiple-point statistics (MPS) simulation method. Such algorithms aim at reproducing not only the 2-point but also higher-order statistics found in a so-called training image (TI). The TI is a gridded 2D or 3D conceptual representation of the target spatial field and can be either continuous or categorical (e.g., geologic facies image). It can either be built from a geologic model or from an observed structure (e.g., outcrop).

Various authors have independently introduced the probabilistic sequential geostatistical resampling (SGR) idea outlined above. The conference paper by Hansen et al. [12] was probably the first to describe the approach, considering so-called block updates where a box-shaped randomly located section of the current model is iteratively resimulated. Almost simultaneously, Fu and Gómez-Hernández [5] proposed a variant of the method that they termed blocking MCMC (BMCMC), which handles multi-Gaussian

\* Corresponding author. Tel.: +32 14333219.

E-mail addresses: [elaloy@sckcen.be](mailto:elaloy@sckcen.be) (E. Laloy), [niklas.linde@unil.ch](mailto:niklas.linde@unil.ch) (N. Linde), [djacques@sckcen.be](mailto:djacques@sckcen.be) (D. Jacques), [gregoire.mariethoz@unil.ch](mailto:gregoire.mariethoz@unil.ch) (G. Mariethoz).

conditional simulation only. Shortly after, Mariethoz et al. [22] presented a SGR algorithm that resimulates a randomly chosen set of pixels/voxels rather than a contiguous block of pixels/voxels. This approach that was named iterative spatial resampling (ISR) was coupled with the direct sampling (DS) MPS algorithm of Mariethoz et al. [23] for resampling of both categorical and continuous priors. Finally, Hansen et al. [13] applied the approach by Hansen et al. [12] to more case studies and clarified the theoretical background of SGR, which they referred to as sequential Gibbs sampling (SGS). Even if all SGR variants presented above fall under the umbrella of Gibbs sampling theory [9], it is worth noting that the latter also forms a common framework for unconditional multi-Gaussian simulation [e.g., [6,20]].

The SGS [12,13] and ISR [22] variants differ only in the geometry of the resimulated grid points, which is a box-shaped area for SGS and a set of points for ISR. For convenience, from here on we will follow Ruggeri et al. [35] and use the generic name “SGR” for both SGS and ISR. To sample from a complex prior, SGR can in principle be implemented with any MPS algorithm. It is however very important that the considered MPS code can condition on a large fraction of grid data points (i.e., resimulating only a small fraction of the model). This is currently achieved only by pixel-based MPS techniques, for example, the DS and SNESIM [38] algorithms. We use DS in this study as it possesses good conditioning capabilities and is memory-efficient and relatively fast.

Ruggeri et al. [35] performed a systematic evaluation of SGR within a multi-Gaussian framework. They compared a gradual deformation [15] proposal mechanism with point and block SGR updates for a synthetic linear geophysical inverse problem using a multi-Gaussian prior and different numbers of measurements and noise levels. Results by Ruggeri et al. [35] suggest that the computational cost of producing one independent realization of the posterior by SGR is often prohibitively large even for relatively simple inverse problems. Ruggeri et al. [35] conclude that this finding warrants further research into model parameter reduction techniques that reduce parameter dimensionality and thus complexity of the inverse problem. This is in line with the work by Laloy et al. [19] who proposed a new reduced multi-Gaussian model parameterization, that is easily coupled with advanced MCMC sampling techniques [e.g., [18,39,41]]. For training-image based inference of non-Gaussian/categorical structures, however, reducing the dimensionality of the parameter field is arguably more difficult. Even though a few model reduction methods have recently been proposed [16,21,40], the conceptual simplicity and flexibility of SGR remain attractive. To the best of our knowledge, no critical analysis of SGR performance has been proposed so far for non-multi-Gaussian cases. Only rather simple problems involving either only 9 data points [22] or unrealistically large measurement errors [13,14] have been considered. Using a very limited number of measurement data and/or large measurement errors makes the likelihood function rather flat and the posterior target is thus easy to sample. It is unclear at this stage whether training-image based SGR can handle more complicated problems with more realistically peaked likelihood functions. Furthermore, the posterior distribution might be multi-modal which, as shown herein, is not easily dealt with by standard SGR.

Parallel tempering (PT) [3,4,10], also called Metropolis-coupled MCMC, consists of parallel Markov chains that sample unnormalized target posterior density functions (pdf) raised to different powers, the inverse of which are called temperatures. The different chains regularly swap their temperatures, with the hot chains sampling a flattened posterior density landscape while the unit temperature(s) chain(s) explore(s) the desired distribution. The hot chains can more easily jump from one basin of attraction of the

posterior to another, and this information is shared through swapping with the cold chain(s) that more intensively explore individual modes. This process can dramatically improve exploration of multi-modal posterior distributions while preserving a theoretically consistent sampling [3,4,10].

Up to now, application of PT to geosciences problems remains limited. In the area of reservoir simulation, Mohamed et al. [26] applied PT to the inversion of the Imperial College fault (ICF) model, considering 3 unknown model parameters and 10 parallel Markov chains. For this application, PT was shown to explore much more efficiently the posterior parameter space than two stochastic optimization algorithms which got stuck within local optima. The study by Carter and White [2] is also focused on posterior exploration of the ICF model, considering from 1 to 13 unknown parameters and using 48 to 64 parallel chains. Carter and White [2] compared a simple random walk Metropolis (RWM) [25] sampler against the same algorithm equipped with PT for an ICF model with one unknown. This clearly demonstrated the superiority of PT for sampling the associated multi-modal posterior parameter distribution. Lastly, Sambridge [36] used as many as 380 parallel chains (and temperatures) to solve a synthetic trans-dimensional (the number of parameters is unknown) geophysical inverse problem for which the true model has 13 unknowns. Results by Sambridge [36] show a spectacular performance improvement by PT in terms of mixing and convergence towards the target data misfit. The use of PT thus appears to be beneficial not only for recovering multi-modal posterior distributions, but also for finding the maximum a posteriori estimate (MAP) of complex unimodal distributions. To date, PT applications in the geosciences have been concerned with rather low-dimensional parameter spaces. We hypothesize that PT may be advantageous for posterior inference in high parameter dimensions as well, such as spatially-distributed subsurface properties. We further suggest that this is possible even when considering a number of levels in the temperature ladder that is very small compared to the dimensionality of the parameter space. For completeness, we note that independently of our work, the idea of coupling PT with SGR also recently appeared as an outlook in the study by Ruggeri et al. [35].

In this paper, we illustrate the limits of the standard SGR method for posterior sampling of categorical fields, and benchmark a PT implementation with respect to both data fitting and diversity of the sampled posterior distribution. We refer to the proposed algorithm as PT-SGR for parallel tempering SGR. In contrast to previous work with PT in the geosciences, the inverse problems considered herein are quite high-dimensional (7575 to 10,000 sampled parameters). Moreover, only a relatively limited amount of parallel chains are used: from 16 to 24. This allows for parallel implementation on workstation computers or small clusters. Besides PT, we also investigate which settings of the SGR algorithm achieve the best performance for categorical field inference. Our proposed PT-SGR approach is demonstrated using synthetic (error corrupted) data from two flow and transport experiments in categorical 10,000- and 7575-dimensional 2D hydraulic conductivity fields that represent a channelized aquifer. These inverse problems involve realistically complex likelihood landscapes, one of which is made bimodal by construction.

This paper is organized as follows. Section 2 presents the different elements of our inversion approach. This is followed in Section 3 with numerical experiments which include a performance analysis of SGR for different algorithmic settings, and a benchmarking of PT-SGR against SGR for the same multi-core computational resources. Section 4 then provides further discussion of the performance and limitations of our method and discusses possible future developments. Finally, Section 5 concludes this paper with a summary of the most important findings.

## 2. Methods

### 2.1. Bayesian inference

A common stochastic representation of the forward problem is

$$F(\boldsymbol{\theta}) = \mathbf{d} + \mathbf{e}, \quad (1)$$

where  $F(\boldsymbol{\theta})$  is a deterministic, error-free forward model that expresses the relation between the unknown parameters  $\boldsymbol{\theta}$  and the measurement data  $\mathbf{d} = (d_1, \dots, d_N) \in \mathbb{R}^N$ ,  $N \geq 1$ , and the noise term  $\mathbf{e}$  lumps all sources of errors.

In the Bayesian paradigm, parameters in  $\boldsymbol{\theta}$  are viewed as random variables with a posterior pdf,  $p(\boldsymbol{\theta}|\mathbf{d})$ , given by

$$p(\boldsymbol{\theta}|\mathbf{d}) = \frac{p(\boldsymbol{\theta})p(\mathbf{d}|\boldsymbol{\theta})}{p(\mathbf{d})} \propto p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{d}), \quad (2)$$

where  $p(\boldsymbol{\theta})$  denotes the prior distribution of  $\boldsymbol{\theta}$  and  $L(\boldsymbol{\theta}|\mathbf{d}) \equiv p(\mathbf{d}|\boldsymbol{\theta})$  signifies the likelihood function of  $\boldsymbol{\theta}$ . The normalization factor  $p(\mathbf{d}) = \int p(\boldsymbol{\theta})p(\mathbf{d}|\boldsymbol{\theta})d\boldsymbol{\theta}$  can be obtained from numerical integration over the parameter space so that  $p(\boldsymbol{\theta}|\mathbf{d})$  is a proper pdf that integrates to unity. The quantity  $p(\mathbf{d})$  is generally difficult to estimate in practice but is not required for parameter inference when the parameter dimensionality is fixed. In the remainder of this paper, we will thus focus on the unnormalized density  $p(\boldsymbol{\theta}|\mathbf{d}) \propto p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{d})$ .

To avoid numerical over- or underflow, it is convenient to work with the logarithm of  $L(\boldsymbol{\theta}|\mathbf{d})$  (log-likelihood),  $\ell(\boldsymbol{\theta}|\mathbf{d})$ , instead of  $L(\boldsymbol{\theta}|\mathbf{d})$ . If we assume  $\mathbf{e}$  to be normally distributed, uncorrelated and with known constant variance,  $\sigma_e^2$ , the component of  $\ell(\boldsymbol{\theta}|\mathbf{d})$  that depends on  $\boldsymbol{\theta}$  can be written as

$$\ell(\boldsymbol{\theta}|\mathbf{d}) = -\frac{1}{2}\sigma_e^{-2} \sum_{i=1}^N [d_i - F_i(\boldsymbol{\theta})]^2, \quad (3)$$

where the  $F_i(\boldsymbol{\theta})$  are the simulated equivalents to the  $i = 1, \dots, N$  measurement data,  $d_i$ .

An exact analytical solution of  $p(\boldsymbol{\theta}|\mathbf{d})$  is not available for the type of inverse problems considered herein. We therefore resort to MCMC simulation to generate samples from the posterior pdf [see, e.g., [31]]. The SGR algorithm independently developed by Hansen et al. [12,13] and Mariethoz et al. [22] is used to approximate the posterior distribution. A detailed description of this sampling scheme can be found in the cited references and a convergence proof is given by Hansen et al. [13]. A brief summary of SGR is given in Section 2.2.

### 2.2. Sequential geostatistical resampling from a training image

For a symmetric proposal distribution, the classical Metropolis acceptance probability,  $\alpha(i, j)$  is given by

$$\alpha(i, j) = 1 \wedge \left( \frac{p(\boldsymbol{\theta}_j)L(\boldsymbol{\theta}_j|\mathbf{d})}{p(\boldsymbol{\theta}_i)L(\boldsymbol{\theta}_i|\mathbf{d})} \right), \quad (4)$$

where the function  $\wedge$  takes the minimum of the left and right hand side numbers. For complex prior models, however, computing  $p(\boldsymbol{\theta})$  might be difficult if not impossible. To overcome this limitation, Mosegaard and Tarantola [27] introduced a different version of the Metropolis algorithm in which the prior probabilities,  $p(\boldsymbol{\theta})$ , need not to be computed. The approach of Mosegaard and Tarantola assumes that a generating algorithm,  $G = q(i, j)$ , exists that is able to sample from  $p(\boldsymbol{\theta})$  directly, in such a way that any proposal,  $\boldsymbol{\theta}_j$ , created by perturbation of the current state,  $\boldsymbol{\theta}_i$ , is itself a draw from  $p(\boldsymbol{\theta})$ . The Metropolis acceptance probability of a move from  $\boldsymbol{\theta}_i$  to  $\boldsymbol{\theta}_j$  can then be reduced to

$$\alpha(i, j) = 1 \wedge \left( \frac{L(\boldsymbol{\theta}_j|\mathbf{d})}{L(\boldsymbol{\theta}_i|\mathbf{d})} \right). \quad (5)$$

Later called “extended Metropolis” sampler [13], this approach forms the basis of SGR. After initializing the chain with  $\boldsymbol{\theta}_i$  drawn from  $p(\boldsymbol{\theta})$ , the latter proceeds in the three following steps to generate a Markov chain. First a candidate model,  $\boldsymbol{\theta}_j$ , is generated by resimulating a random fraction of the current state,  $\boldsymbol{\theta}_i$ , according to the prior model distribution. Since the conditioning points are chosen at random, this mechanism corresponds to a symmetric proposal distribution,  $q(i, j)$ , thus honoring the detailed balance condition:  $q(i, j) = q(j, i)$  [see, e.g. [31], for theoretical details about MCMC]. Next,  $\boldsymbol{\theta}_j$  is either accepted or rejected using Eq. (5). Finally, the chain either moves to  $\boldsymbol{\theta}_j$  ( $\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_j$ ) if the proposal is accepted, or remains at its current location ( $\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i$ ) otherwise. Upon convergence of the chain, the generated model states constitute a set of representative draws from the posterior pdf.

Priors with complex structures can be handled by using a MPS algorithm that samples from a prescribed training image, which acts as prior model. Similarly as Mariethoz et al. [22], we use herein the DS method by Mariethoz et al. [23] as generating and conditioning algorithm. The selected TI is the most classical  $250 \times 250$  binary image representing a channelized aquifer (not shown) that was introduced by Strebel [38] [see also Fig. 4a in [23]].

### 2.3. Parallel tempering sequential geostatistical resampling

In parallel tempering [3,4,10], a temperature ladder,  $\mathbf{T} = [T_1, \dots, T_n]$  with  $T_1 = 1 < T_2 < \dots < T_n$ , is used to increasingly flatten either the posterior density

$$p(\boldsymbol{\theta}, T|\mathbf{d}) \propto [p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{d})]^{1/T}, \quad (6)$$

or the likelihood function

$$p(\boldsymbol{\theta}, T|\mathbf{d}) \propto p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{d})^{1/T}, \quad (7)$$

using a so called temperature  $T$  [11]. When  $T \rightarrow \infty$  in Eq. (7), the tempered distribution  $p(\boldsymbol{\theta}, T|\mathbf{d})$  becomes the prior distribution  $p(\boldsymbol{\theta})$ . Conversely,  $p(\boldsymbol{\theta}, T|\mathbf{d})$  becomes the posterior pdf,  $p(\boldsymbol{\theta}|\mathbf{d})$ , when  $T = 1$ . In this work,  $p(\boldsymbol{\theta})$  does not need to be calculated explicitly as the MPS algorithm generates proposals directly from  $p(\boldsymbol{\theta})$  which is formed by the TI. In this case, Eqs. (6) and (7) reduce to  $p(\boldsymbol{\theta}, T|\mathbf{d}) \propto L(\boldsymbol{\theta}|\mathbf{d})^{1/T}$ .

Each tempered chain undergoes two possible moves throughout sampling: within-chain and between-chain proposals. In our proposed PT-SGR implementation, the within-chain proposal consists of a standard SGR update where a random fraction of the current state of the chain is resimulated according to the TI. The between-chain proposal consists of a swap of model states at two temperature levels  $i$  and  $j$

$$[(\boldsymbol{\theta}_i, T_i), (\boldsymbol{\theta}_j, T_j)] \rightarrow [(\boldsymbol{\theta}_i, T_j), (\boldsymbol{\theta}_j, T_i)], \quad (8)$$

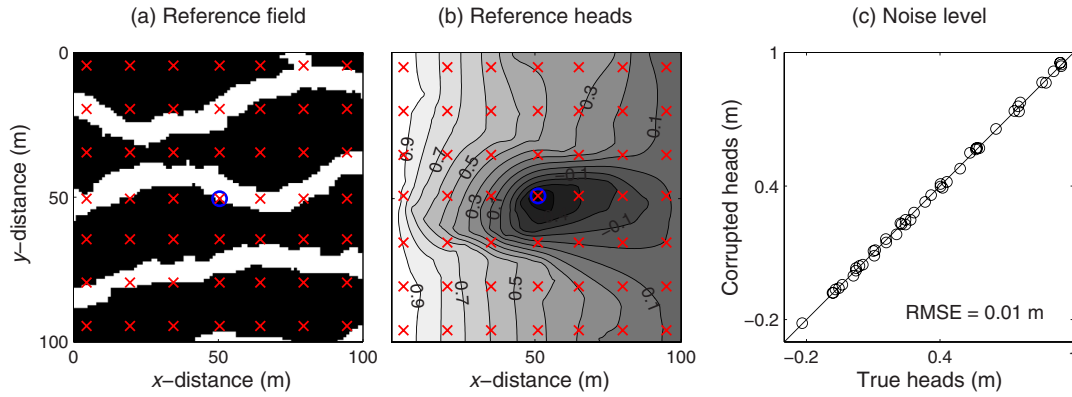
where  $\boldsymbol{\theta}_i$  and  $\boldsymbol{\theta}_j$  are the model parameter vectors in chains  $i$  and  $j$  immediately before the proposed swap. Exchange swap proposals improve the sampling at two levels. At the beginning of the search, they make it easier for the unit temperature ( $T = 1$ ) chain to access regions of the model space with high posterior probabilities that are well removed from its current position. After burn-in, they allow for the unit temperature chain to jump between multiple peaks of the posterior density landscape.

Using Eq. (7), the Metropolis acceptance probability,  $\alpha_s(i, j)$ , of an exchange swap between models  $\boldsymbol{\theta}_i$  and  $\boldsymbol{\theta}_j$  with temperatures  $T_i$  and  $T_j$  is given by Geyer [11]

$$\alpha_s(i, j) = 1 \wedge \frac{L(\boldsymbol{\theta}_j|\mathbf{d})^{1/T_i} p(\boldsymbol{\theta}_j)L(\boldsymbol{\theta}_i|\mathbf{d})^{1/T_j} p(\boldsymbol{\theta}_i)}{L(\boldsymbol{\theta}_j|\mathbf{d})^{1/T_j} p(\boldsymbol{\theta}_j)L(\boldsymbol{\theta}_i|\mathbf{d})^{1/T_i} p(\boldsymbol{\theta}_i)}. \quad (9)$$

Canceling the equivalent terms in the numerator and denominator and performing some reorganization leads to

$$\alpha_s(i, j) = 1 \wedge \left[ \frac{L(\boldsymbol{\theta}_j|\mathbf{d})}{L(\boldsymbol{\theta}_i|\mathbf{d})} \right]^{(1/T_i - 1/T_j)}, \quad (10)$$



**Fig. 1.** (a) Reference categorical field, (b) associated heads and (c) noise-corrupted measurement data used for case study 1. In subfigures a and b, the blue circle marks the location of the pumping well and the red crosses indicate piezometers. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

where the prior distributions  $p(\theta_i)$  and  $p(\theta_j)$  do not appear, thereby allowing us to couple PT with SGR.

The temperature swapping (Eq. (8)) is often restricted to neighboring temperatures [3], either by considering all pairs of neighbors at once [e.g., [29]] or only one pair at the time [e.g., [37]]. Other authors [e.g., [36]] instead proposed to randomly swap models independently of their temperature levels. All these schemes are valid in the sense that the unit temperature chain(s) will ultimately converge to the proper stationary distribution, provided that any given temperature is involved in possible exchange with no more than one other temperature [29] at each proposal step.

A general guideline is that exchange swaps must not happen too frequently such that (1) swapping-induced correlation in the tempered chains is reduced and (2) the risk for the unit temperature chain(s) to get trapped by cycling locally within a certain temperature interval is minimized [28–30]. Nonetheless, it has also been shown that optimal sampling performance is attained with a relatively high frequency of the exchange swap proposals [28]. The mean acceptance probability,  $\alpha_s$ , of a swap move for a given pair (Eq. (10)) is recognized to be another important diagnostic of parallel tempering performance. Obviously, an excessively small  $\alpha_s$  will hamper exploration by the unit temperature chain. As  $\alpha_s$  increases, however, the tempered chains will tend to keep exchanging each other's models without creating new configurations, thereby slowing down posterior sampling [28]. Optimal  $\alpha_s$  values of about 20% [30] and 39% [29] have been proposed under certain conditions, whereas good sampling performance was found with a  $\alpha_s$  value as low as 8% [28].

A pseudo-code of the proposed PT-SGR algorithm is given in Algorithm 1.

As discussed above, there are several options for selecting temperature pairs at swapping time. Two such options are described in Algorithm 2: considering either randomly located or adjacent temperatures, with each temperature level involved in one exchange swap. In the case of adjacent temperatures, the same pairs can of course not be selected every time. When swapping occurs, the chosen pairs are thus either (1, 2), (3, 4), ... or (2, 3), (4, 5), ... depending on whether the iteration number is even or odd [e.g., [29]].

### 3. Case studies

#### 3.1. Case study 1: steady-state flow

Our first synthetic case study considers steady-state head data collected at various locations within a channelized 2D aquifer (Fig. 1a). The 100 × 100 modeling domain lies in the  $x$ – $y$  plane

with a grid cell size of 1 m. Channels and matrix are assigned hydraulic conductivity values of  $1 \times 10^{-2}$  m/s and  $1 \times 10^{-4}$  m/s, respectively. Steady state groundwater flow is simulated using MaFloT [17] which is a finite-volume algorithm for 2D flow and transport in porous media. We assume no flow boundaries at the upper and lower sides and fixed head boundaries on the left and right sides of the domain so that a lateral head gradient of 0.025 (–) is imposed, with water flow in the  $x$ -direction. A pumping well extracting  $0.003 \text{ m}^2/\text{s}$  is located at the center of the domain. Simulated heads are collected at 49 locations that are regularly spread over the domain (Fig. 1a and b). These data were then corrupted with a Gaussian white noise using a standard deviation of 0.01 m, leading to a root-mean-square-error (RMSE) of 0.01 m for the measurement data (Fig. 1c). This translates into a reference log-likelihood (component that depends on  $\theta$ ),  $\ell(\theta|\mathbf{d})$ , of –24.5.

#### 3.2. DS settings for case study 1

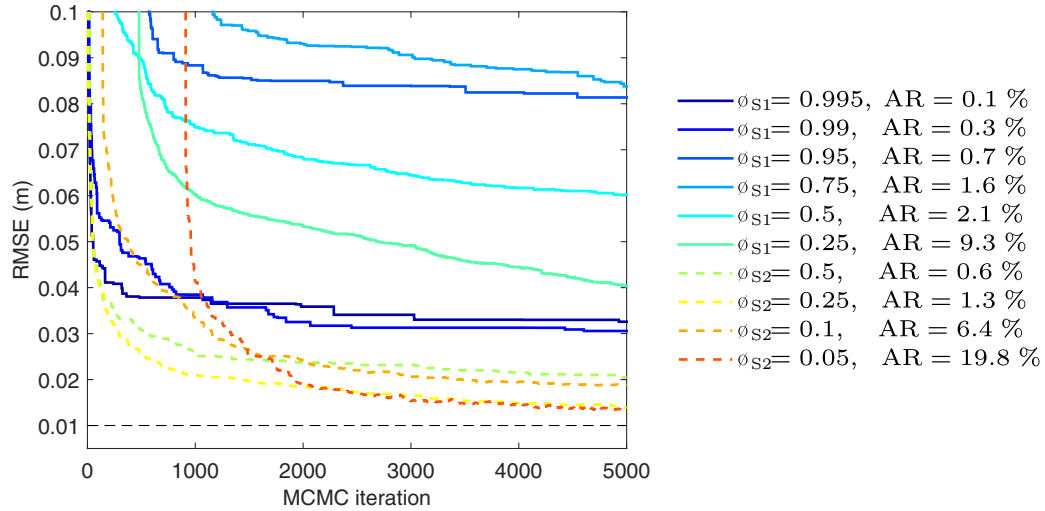
The parameters of the DS simulation used for case study 1 are a neighborhood made of 50 nodes and a distance threshold set to 0.05. This means that for any simulated node, the data event (pattern) made of the 50 closest neighbors is considered, and up to only 2 mismatching nodes are allowed [see [23], for details]. The maximum scanned fraction of the TI is set to 0.9.

#### 3.3. SGR settings

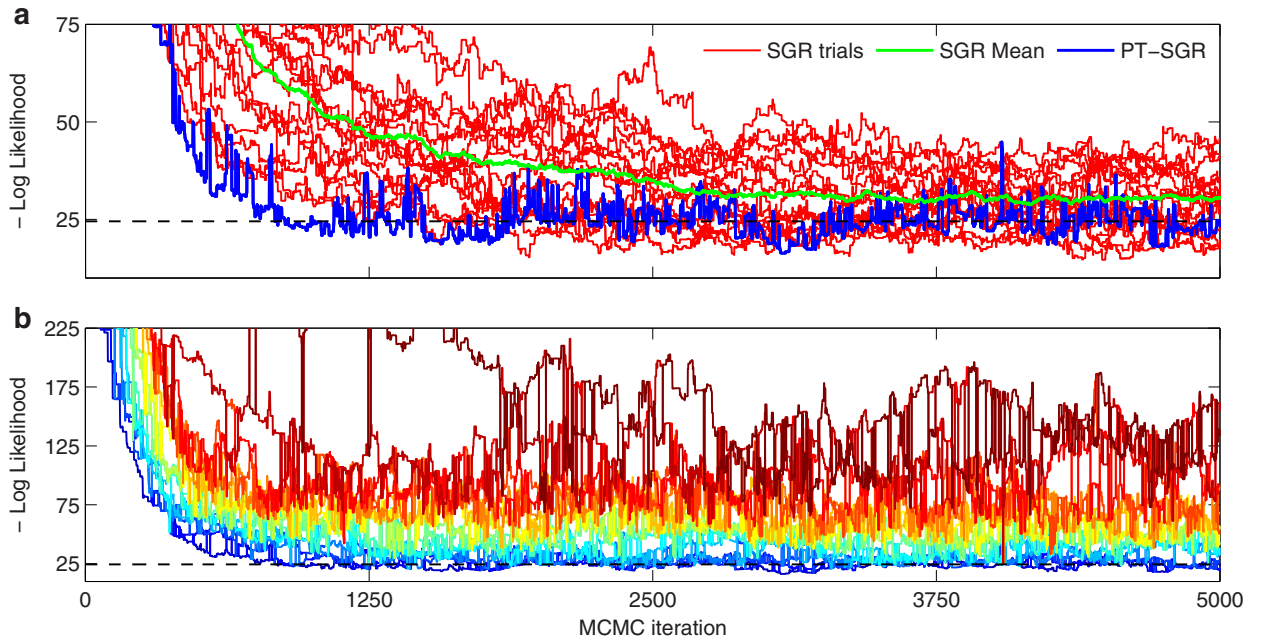
Apart from the employed MPS algorithm, important SGR algorithmic settings are (1) the type of conditioning, that is, whether the pixels to resimulate are defined by a set of points that are distributed throughout the model domain or if they all belong to a box-shaped area, and (2) the size of the randomly located model fraction that is resimulated,  $\phi$ . The latter can be fixed beforehand, adapted during burn-in of the MCMC sampling or drawn randomly from a certain probability distribution. All of these three options are explored in this study. In this section, we study the sampling performance achieved by conditioning on points (S1), or on all the points outside a square block (S2) for different sizes of  $\phi$ . For S1, the six following fractions were considered: 0.995, 0.99, 0.95, 0.75, 0.5 and 0.25. As of S2,  $\phi$  was set to 0.5, 0.25, 0.1 and 0.05. For each combination of settings, the setup of case study 1 (Section 3.1) was used to perform 4 different MCMC trials for a total of 5000 forward model runs. In this and all other MCMC experiments conducted in this study, we initialized each Markov chain by randomly sampling  $p(\theta)$ .

Fig. 2 displays the resulting sampled root-mean-square error (RMSE) trajectories and mean acceptance rate (AR) of the MCMC.





**Fig. 2.** Trace plot of the mean sampled RMSE values across 4 repetitions for the tested conditioning strategies. Solid and dashed colored lines denote resimulating a set of points (S1) and a box-shaped area (S2), respectively. Each color represents a given size of the (randomly located) model fraction that is resimulated. The dashed black line signifies the true RMSE. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).



**Fig. 3.** (a) Trace plot of the sampled negative log-likelihood values by the unit temperature chain evolved by PT-SGR (blue line) and the 16 independent SGR trials (red lines) for case study 1. The green line denotes the mean trajectory of the 16 SGR trials. (b) Trace plot of the sampled negative log-likelihood values by the 16 PT-SGR chains with each temperature coded with a different color. The temperature increases as the color varies from dark blue (temperature index of 1) to dark red (temperature index of 6). In both subfigures, the horizontal dashed black line denotes the true negative log-likelihood of 24.5, corresponding to a RMSE of 0.01 m. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

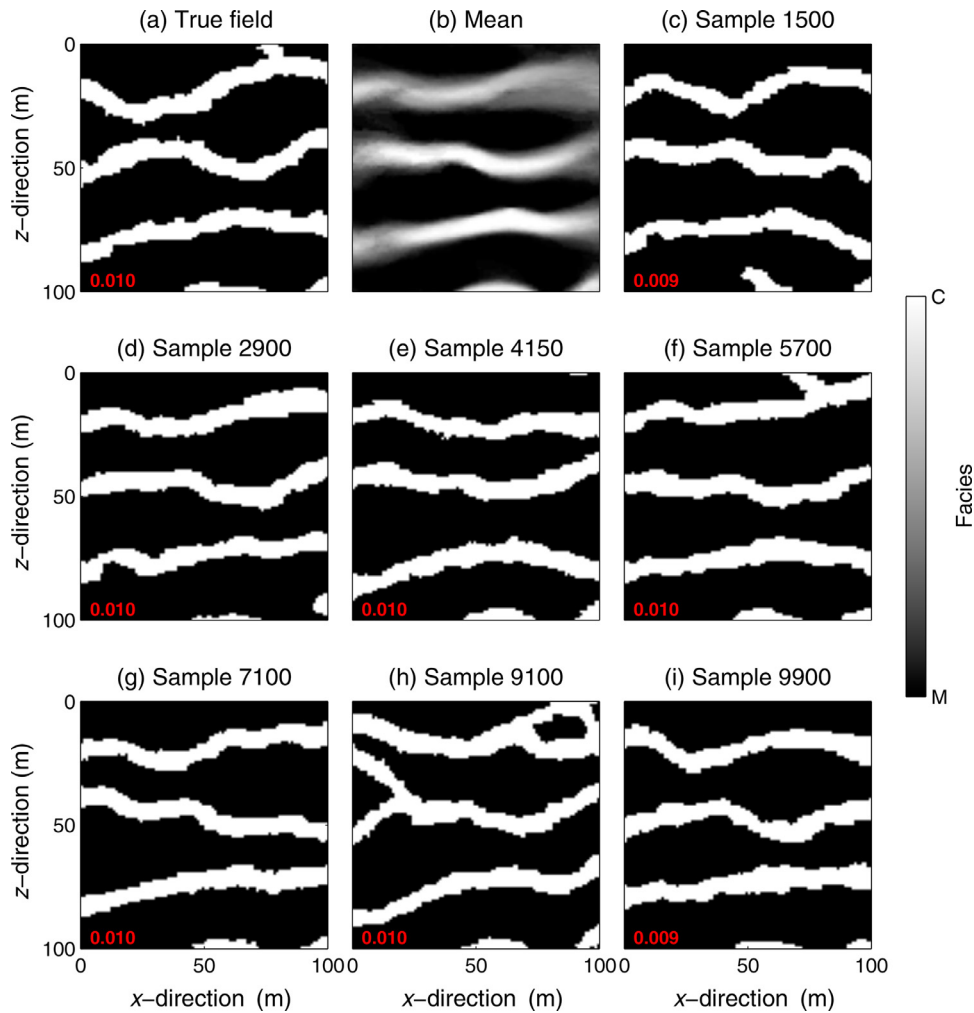
Averages of the 4 trials are presented. Clearly, resimulating a box-shaped area (S2) shows a superior performance with respect to data fitting. As expected, the AR decreases with  $\phi$  for both S1 and S2. Large resimulated fractions induce low AR values, below 1 or 2% ( $\phi_{S1} = 0.995$  to  $\phi_{S1} = 0.5$  and  $\phi_{S2} = 0.5$  and  $\phi_{S2} = 0.25$ ). Such small AR values characterize a prohibitively slow evolution of the MCMC chain.

Based on the above findings, we decided to use the resimulation strategy S2 in all of the following tests. Since the optimal value of  $\phi$  is likely to depend on the problem at hand, in the remainder of this paper and unless stated otherwise  $\phi_{S2}$  is tuned online to try to reach an AR value of 20% during the first 10% of the MCMC iterations. The motivation for this target value is based on the fact that an AR of about 23% is considered optimal for Gaus-

sian proposal and target distributions whereas an AR in the range 10%–50% is generally recommended [32,33].

### 3.4. Parallel tempering settings

For the case studies considered in this paper, limited testing with the different swapping strategies described in Section 2.3 showed no overwhelming advantage of any specific strategy. Nevertheless, randomly proposing to swap model states after every regular within-chain MCMC update appeared to be the most robust and efficient approach. We therefore do so for all of our numerical experiments. With respect to the  $\alpha_s$  values of the individual tempered chains, we use a common loglinear temperature ladder [2,34] with maximum level such that  $\alpha_s$  is (almost) always com-



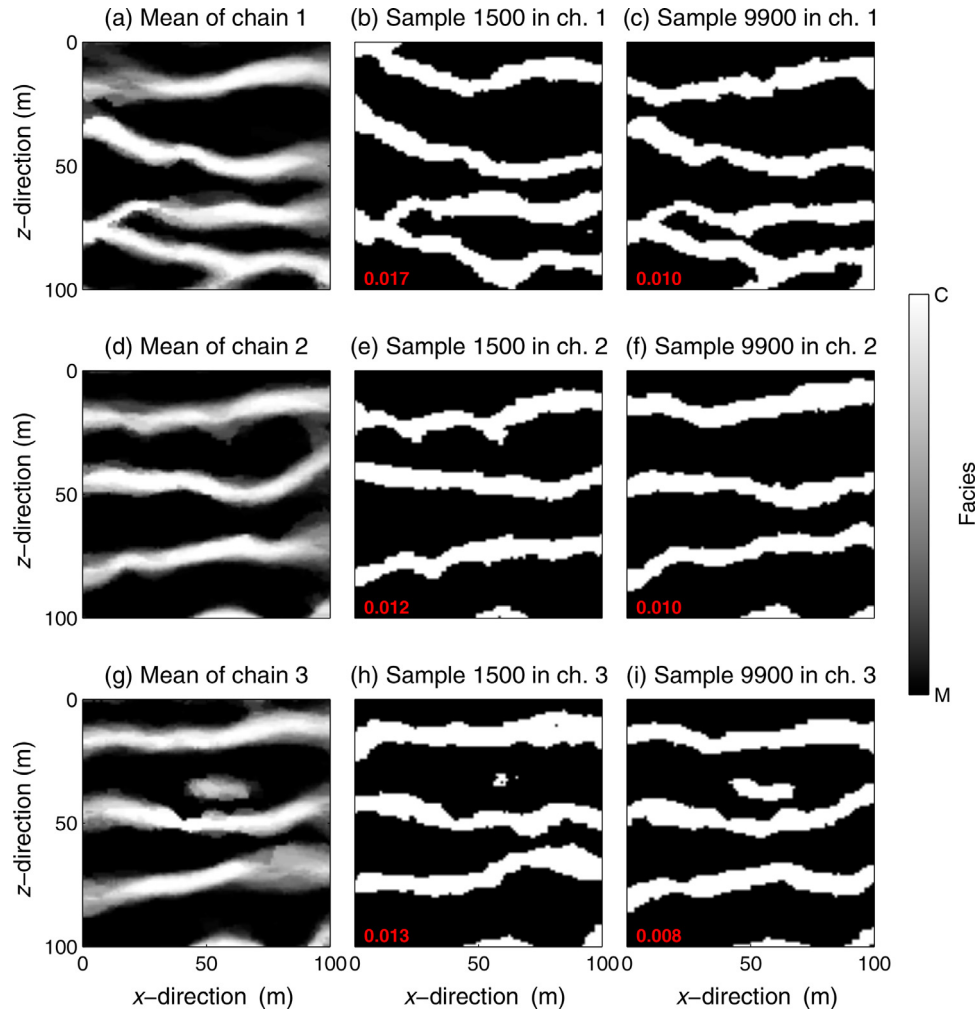
**Fig. 4.** Posterior mean and 8 successive posterior realizations taken at regular intervals throughout sampling for the PT-SGR trial of case study 1. The posterior mean is computed on the basis of the samples produced by the unit temperature chains after a burn-in of 1500 MCMC iterations and using a thinning factor of 50, thus leading to a total of 170 posterior samples. The red number in the lower left corner of each plot is the corresponding RMSE (m). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

prised between 5% and 30%. The rationale behind a loglinear scale is that a pair of neighboring cool chains ( $T = 1$  or slightly higher) likely needs a smaller temperature difference for an exchange swap to be accepted compared to a pair of hotter chains. Other settings that perform better may very well exist, the quest for which is beyond the scope of this study.

### 3.5. Convergence of the Markov chain Monte Carlo simulation

The use of multiple (unit temperature) Markov chains makes it possible to use of the potential scale reduction factor,  $\hat{R}$  [8], for monitoring convergence of the MCMC sampling. The  $\hat{R}$  statistic compares for each parameter of interest the average within-chain variance to the variance of all the chains mixed together. The closer the values of these two variances, the closer to unity the value of  $\hat{R}$ . Values of  $\hat{R}$  smaller than 1.2 are commonly deemed to indicate convergence to a limiting distribution. In principle,  $\hat{R}$  offers a stronger convergence assessment than merely considering the moment when the sampled (log-)likelihood (and thus RMSE) values reach an equilibrium. The latter indeed signifies only that the posterior distribution has been located, whereas the former aims at evaluating whether it has been adequately explored. For

example, in the study by Laloy et al. [19] it was found that 25 times more MCMC steps were needed to appropriately explore a 1000-dimensional posterior than to start sampling it. For the considered case studies and computational budgets, our simulation results indicate that  $\hat{R}$ -convergence is never achieved by SGR, no matter whether computed on the basis of 3 (randomly chosen) chains or all of the 16 (case study 1) or 24 (case study 2) independent trials. Unfortunately, the  $\hat{R}$  statistic is not well suited for monitoring convergence of several unit temperature chains within a given tempered ensemble. Indeed, running PT-SGR for case study 1 with 3 out of the 16 temperatures set to 1 results in an exaggeratedly fast  $\hat{R}$ -convergence (not shown). The 10,000 individual  $\hat{R}$  values may even jointly fall below 1.2 after less iterations than required for the unit chains to sample the appropriate likelihood values. The reason for this is that the swapping dynamics causes large correlations between states/models of neighboring cool chains. The within-chain variances thus become similar enough for  $\hat{R}$  to be satisfied prematurely. We therefore refrain from using the  $\hat{R}$  diagnostic. Instead we simply resort to the point in time when (log-)likelihood values start to fluctuate around a constant level to define burn-in. From this moment on our algorithm starts drawing posterior samples. One must bear in mind, however, that given the



**Fig. 5.** Mean sampled model over MCMC iterations 1500–10,000 (using a thinning factor of 50), and sampled models after 1500 and 9900 MCMC iterations for 3 out of the 16 independent SGR chains and case study 1. The red number in the lower left corner of each plot is the corresponding RMSE (m). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

large dimensionality of the considered problems, it is evident that the posterior target is not fully explored within our limited computational budget (10,000 to 25,000 MCMC iterations) and we do not claim to do so. By a convenient abuse of terminology, we nevertheless refer to the resulting set of posterior samples as the “posterior distribution”.

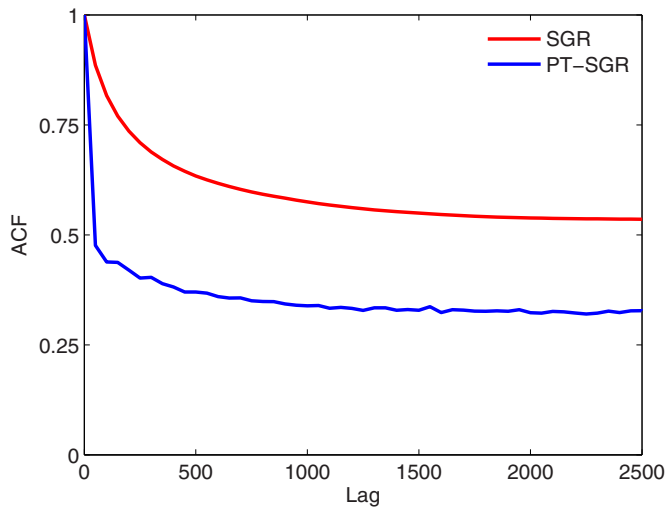
### 3.6. Inversion results for case study 1

For this case study, a total of 10,000 MCMC iterations is allowed for both SGR and PT-SGR. For a classical single-chain SGR trial, this translates into 10,000 forward model evaluations. The PT-SGR algorithm with  $n$  temperature levels is best run with parallel calculation of both the  $n$  DS simulations and  $n$  forward model evaluations performed per MCMC iteration. This roughly leads to a similar CPU-time per MCMC iteration (from 5 to 10 s herein) between SGR and PT-SGR. One must note, however, that some minor additional computational time is needed for PT-SGR due to communication overhead, the extent of which depends on hardware- and software-specific details. Here PT-SGR is ran on a multi-core platform, with  $n = 16$ . A loglinear temperature ladder was selected between unity and a maximum temperature of 6, together with a single unit temperature chain.

Fig. 3a depicts the sampled negative log-likelihood ( $-\ell(\theta|\mathbf{d})$ ) trajectories for the unit temperature PT-SGR chain and the 16 in-

dependent SGR trials. It is observed that the PT-SGR trial samples appropriate mean RMSE values after some 1250 iterations. In contrast, the basic SGR algorithm shows a large spread of trajectories. Overall, the PT-SGR chain converges towards the reference data misfit at least as fast as the fastest of the 16 SGR chains. It also takes about 8000 MCMC iterations for the mean of the 16 SGR trials to reach the target RMSE value (not shown). For this particular run, this leads to a 6 times speed-up of PT-SGR. Limited additional testing with PT-SGR confirmed (I) a similar data fitting efficiency of PT-SGR to that of the best performing SGR trial and (II) a 5–8 times speedup of PT-SGR for locating the posterior compared to the mean SGR behavior. This speedup is accomplished by the (random) mixing across the whole temperature ladder (Fig. 3b). Across the tempered PT-SGR chains, the AR associated with the regular and swap moves,  $\alpha$  and  $\alpha_s$ , are 24% (with range of 19%–26%) and 19% (with range of 6%–24%), respectively.

The posterior distribution sampled by PT-SGR is illustrated in Fig. 4 that shows the reference field together with the posterior mean and 7 successive posterior realizations from the unit temperature chain. The posterior mean in Fig. 4 resembles the true model (Fig. 1a) relatively well and the derived posterior uncertainty is rather small (compare realizations c–i in Fig. 4). With respect to SGR, each of the 16 sampling runs turns out to remain in a specific region of the model space, as depicted by Fig. 5. Indeed, the variability within the individual Markov chains is quite limited:



**Fig. 6.** Mean autocorrelation function (ACF) of the 10,000 conductivity grid values derived from PT-SGR (blue line) or SGR (red line) for lags 0–2500 and case study 1. The lag- $k$  autocorrelation is defined as the correlation between draws  $k$  lags apart. Listed statistics are computed for the last 8500 iterations of the unit temperature chain of PT-SGR or the 16 independent SGR chains, using a thinning factor of 50 thereby leading to a set of 170 sampled models for each chain. For SGR, the average of the 16 chains is presented. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

models sampled more than 8500 iterations apart look very similar both to each other and to the chain average. Though also limited, the variability sampled by the (unit temperature) PT-SGR chain is nevertheless larger than for any given SGR trial. This is confirmed by the mean autocorrelation functions (ACF) calculated for the two tested algorithms (Fig. 6): the ACF of PT-SGR drops much more rapidly than that of SGR, and stabilizes around a 1.6 times smaller value: 0.33 against 0.54 for SGR. The stabilization around a value larger than zero is caused by the fact that some specific binary grid elements never change of value throughout the considered set of MCMC draws. Label switching for these grid elements is proposed but the resulting models have always too low likelihood to be accepted by the Markov chain. An ACF value of 1 is thus assigned to these grid elements which influences the mean ACF. Herein 33% of the 10,000 grid elements have not been updated for PT-SGR, against up to 54% for SGR.

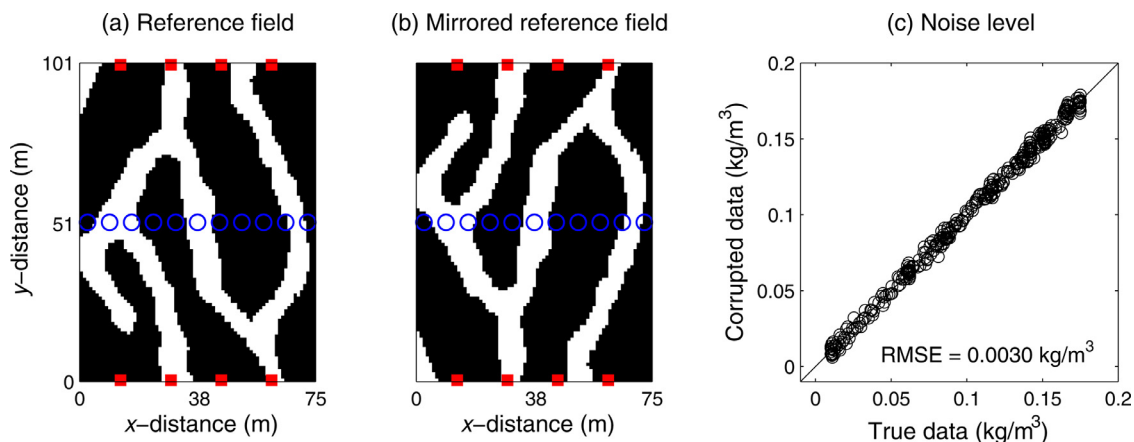
The above results show that for case study 1, running PT-SGR with 16 parallel chains is a better option than running 16 in-

dependent SGR chains. For the CPU budget needed by PT-SGR to start sampling the posterior, most of the SGR chains are still exploring parts of the prior that do not belong to the posterior (Fig. 3a). Indeed, it requires about 6 times more computational time for the 16 SGR chains to jointly sample the posterior. As of posterior diversity, each SGR trial gets trapped in a small region of the posterior model space (Fig. 5). The situation is arguably better for PT-SGR (see Figs. 4 and 6) even though it is far from having explored the full posterior range.

### 3.7. Case study 2: tracer experiment

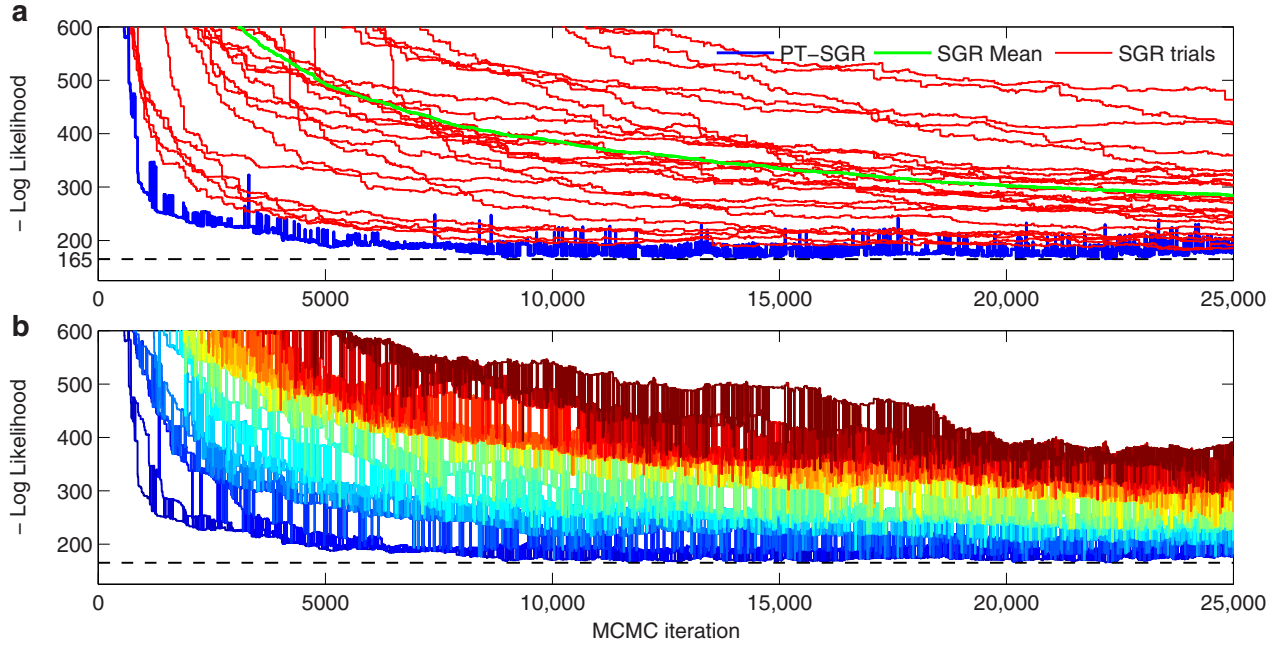
Our second case study uses simulated tracer breakthrough curves at different wells as measurement data. The modeling domain is  $75 \times 101$  and is located in the  $x-y$  plane with a grid cell size of 1 m. Channels and matrix are again assigned hydraulic conductivity values of  $1 \times 10^{-2}$  m/s and  $1 \times 10^{-4}$  m/s, respectively. Steady state groundwater flow and conservative transport are both simulated using MaFloT. No-flow boundaries are assumed at the left and right sides, and fixed head boundaries on the upper and lower sides of the domain. These fixed heads are set to 0 m at both sides, and 11 pumping wells individually extracting  $0.0005 \text{ m}^2/\text{s}$  of water are spaced 7 m apart along the horizontal line located at equal distance from the top and bottom sides (Fig. 7a). The facies values at the 11 wells are assumed to be known exactly and serve for direct conditioning. A conservative tracer with concentration of  $1 \text{ kg m}^{-3}$  is applied within 8 model cells of the top and bottom boundaries using a step function. The  $x-y$  coordinates of these cells are (14,1), (30,1), (46,1), (62,1), (14,101), (30,101), (46,101), and (62,101) (Fig. 7a). The background solute concentration is assumed to be  $0.01 \text{ kg m}^{-3}$ . Ignoring density effects, conservative transport of the tracer through the subsurface is simulated using open boundaries on all sides, and longitudinal and transverse dispersivities both set to 0.1 m. Solute transport was monitored during a period of 10 days with concentration measurements made every 8 h in the 11 extraction wells, resulting into a total of 330 observations. These simulated data were then corrupted with a Gaussian white noise using a standard deviation equivalent to 3% of the mean observed concentration. This led to root-mean-square-error (RMSE) and log-likelihood ( $\ell(\theta|\mathbf{d})$ ) of  $0.0030 \text{ kg m}^{-3}$  and  $-165$ , respectively, for the measurement data (Fig. 7c).

This setup has the attractive feature of causing the posterior facies distribution to include two distinct modes with equal probability. Indeed, the reference field of Fig. 7a and its mirrored image



**Fig. 7.** (a) Reference categorical field, (b) associated symmetric (mirrored) field and (c) noise-corrupted measurement data used for our second synthetic case study. In subfigures a and b, the red squares denote the application points of the tracer and the blue circles mark the locations of the pumping wells where concentrations are monitored. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).





**Fig. 8.** (a) Trace plot of the sampled negative log-likelihood values by the unit temperature chain evolved by PT-SGR (blue line) and the 24 independent SGR trials (red lines) for case study 2. The green line denotes the mean trajectory of the 24 SGR trials. (b) Trace plot of the sampled negative log-likelihood values by the 24 PT-SGR chains with each temperature coded with a different color. The temperature increases as the color varies from dark blue (temperature index of 1) to dark red (temperature index of 2). In both subfigures, the horizontal dashed black line denotes the true negative log-likelihood of 165, corresponding to a RMSE of 0.003 kg/m<sup>3</sup>. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

shown in Fig. 7b both lead to the exact same simulated concentration data and thus likelihood function value. We thus consider this rather challenging case study to be especially instructive as the posterior target is known to present (at least) two separate modes.

### 3.8. DS settings for case study 2

The parameters of the DS simulation used for case study 2 are a neighborhood made of 75 nodes, a distance threshold of 0.01 and a maximum scanned fraction of the TI of 0.9.

### 3.9. Inversion results for case study 2

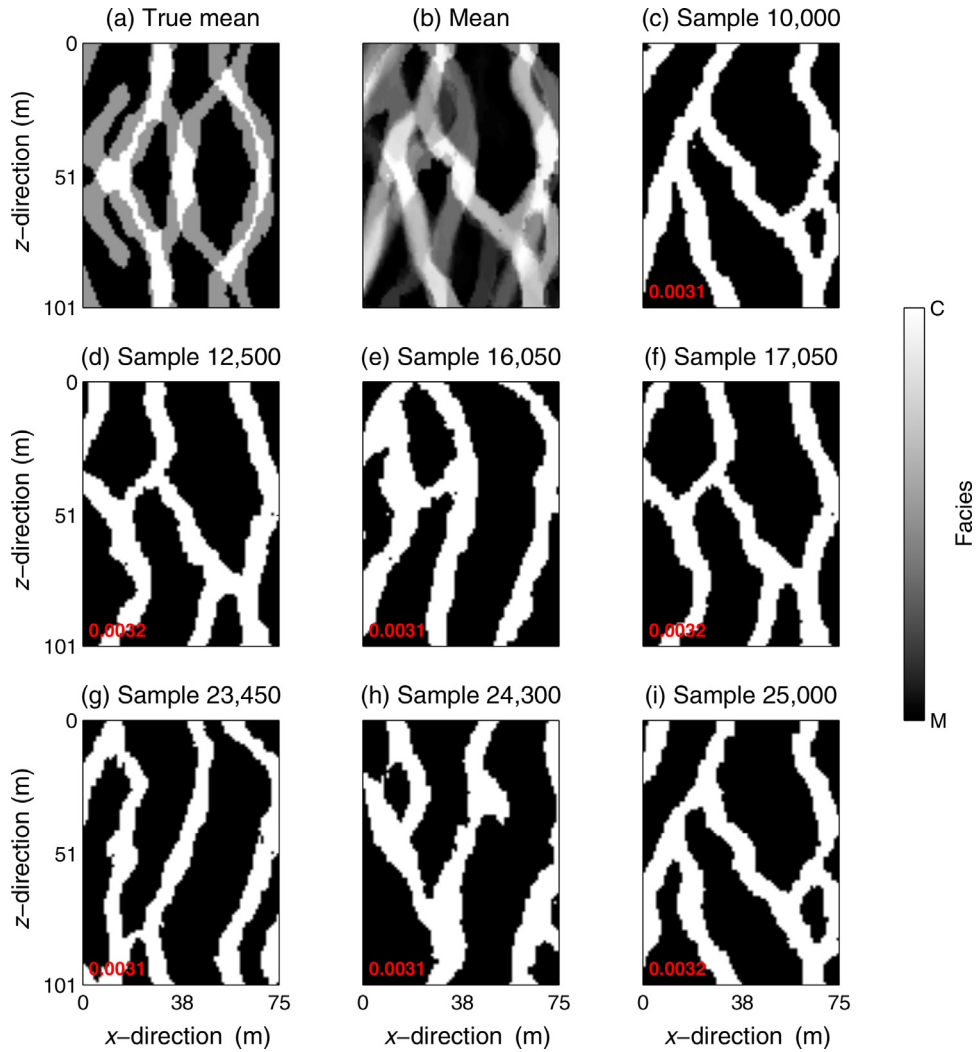
A total of 25,000 MCMC iterations is used for this case study. The PT-SGR algorithm is again run on a multi-core machine with  $n = 24$  and the computational cost incurred by 1 MCMC iteration is in the range of 20–30 s for this workstation. A loglinear temperature ladder between unity and a maximum temperature of 2 is selected together with a single unit temperature chain. Using such a small maximum temperature was needed given the trade-off between the number of available parallel cores and the complexity of the (log-)likelihood landscape. For instance, the peakier the likelihood function, the smaller the temperature intervals need to be for  $\alpha_s$  to be significantly larger than zero. With 24 temperature levels, using a larger maximum temperature than 2 essentially results in a frequency of accepted swaps that is impractically low. In addition, the update mechanism of  $\phi$  that is described in Section 3.3 was slightly modified in an attempt to generate more diverse proposals while keeping the acceptance rate of regular MCMC moves reasonably high. Rather than a tuned constant resimulation block size,  $\phi$  was taken as the (adapted) mean of a triangular pdf bounded between 0.01 and 0.25 from which the actual size of the block to be resimulated was drawn. This was similarly done for both SGR and PT-SGR.

The sampled  $-\ell(\theta|\mathbf{d})$  trajectories for the unit temperature PT-SGR chain and the 24 independent SGR trials are presented in Fig. 8a. The PT-SGR unit chain evolves towards the reference

$-\ell(\theta|\mathbf{d})$  value much faster than the fastest of the SGR trials. Furthermore, the spread of the SGR trajectories is rather large. After 25,000 iterations, only 4 trials are sampling  $-\ell(\theta|\mathbf{d})$  values in the range sampled by PT-SGR and 5 trials are still exploring areas associated with twice as large  $-\ell(\theta|\mathbf{d})$  values. Eq. (5) can be used to calculate the probability of a direct jump from the reference model to the most likely model found by PT-SGR on the one hand, and the most likely model among the 24 SGR chains on the other hand. Doing so reveals that the most likely model sampled by PT-SGR is more than  $1 \times 10^7$  times more likely than that of SGR. The PT-SGR algorithm thus clearly outperforms SGR for this case study.

Even if PT-SGR surpasses SGR, it is evident that the unit temperature PT-SGR chain fluctuates around a slightly larger  $-\ell(\theta|\mathbf{d})$  value than the reference value of  $-165$ . In fact, for iterations 10,000–25,000 the mean sampled  $-\ell(\theta|\mathbf{d})$  value by PT-SGR exceeds the reference value of 165 by 9% and the sampled range actually never contains it (Fig. 8a). This means that the samples produced by PT-SGR are not representative of the posterior distribution. That said, this inverse problem is much more difficult to solve than for case study 1 (see Section 3.6). This is because (I) the large amount of good quality (moderately corrupted) measurements (330) causes the two well separated (log-) likelihood modes to be more peaky, and (II) using transport data induces a more nonlinear relationship between model (parameters) and (log-) likelihood than using steady-state head data.

The AR associated with the PT-SGR chains are lower than for case study 1 but still acceptable: across the whole temperature ladder,  $\alpha$  and  $\alpha_s$ , are 8% (with range of 7%–9%) and 14% (with range of 6%–20%), respectively. The corresponding swap exchange dynamics looks visually good from iteration 10,000 onwards (Fig. 8b). The PT-SGR unit temperature chain does however not mix well. Indeed, the chain basically cycles over the (nearly) same 6–7 models during the last 15,000 MCMC iterations (Fig. 9). The reason for this is likely twofold. First, the maximum temperature of 2 does not flatten the likelihood enough for sufficient exploration by the hot chains. Second and most important, for this rather complex



**Fig. 9.** Mean and 8 successive model realizations taken at regular intervals throughout sampling for the PT-SGR trial of case study 2. The sample mean is computed on the basis of the samples produced by the unit temperature chains over iterations 10,000–25,000 and using a thinning factor of 50, thus leading to a total of 300 samples. The red number in the lower left corner of each plot is the corresponding RMSE (m). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

likelihood landscape the exchange swaps appear to be mostly performed in cycles between a few neighboring temperature levels with colder (hotter) samples almost never traveling to the highest (lowest) levels (not shown). The only way to solve this problem would therefore be to use a ladder with a much larger number of levels and a wider temperature range.

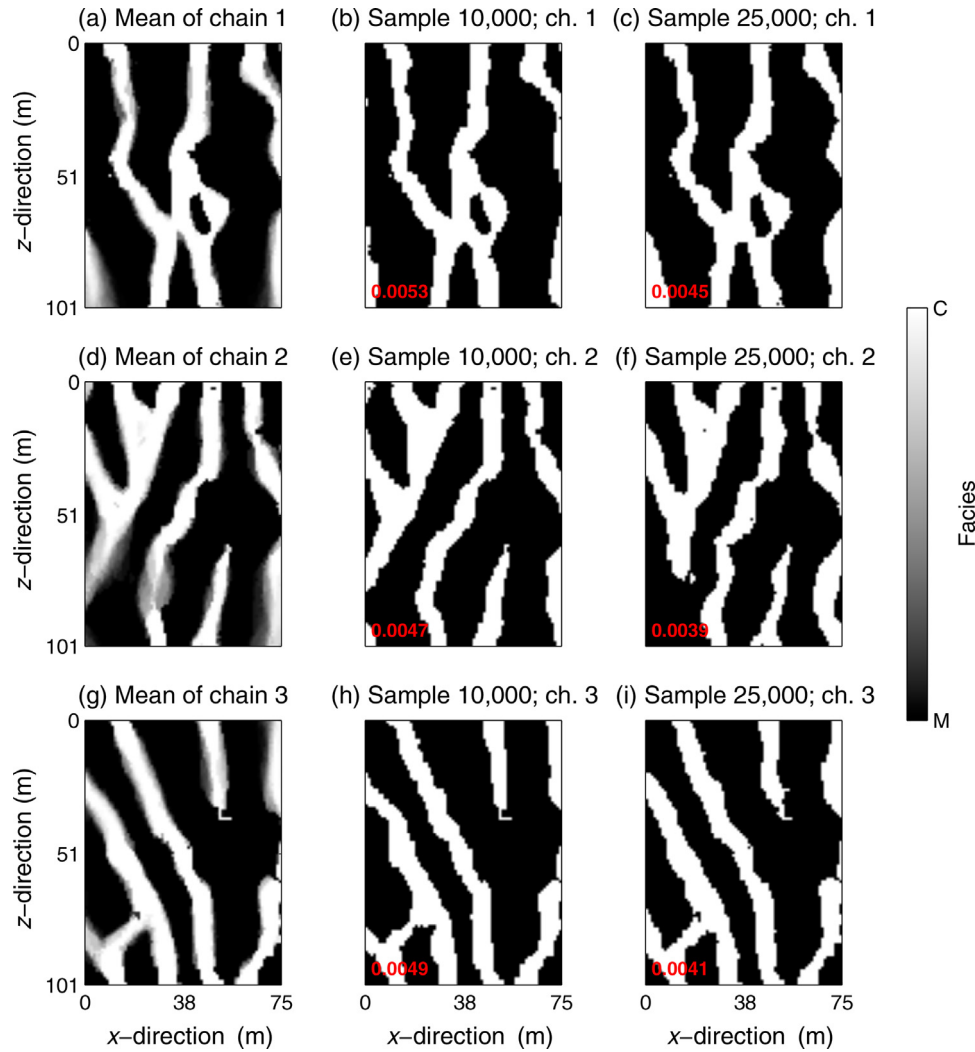
As depicted in Fig. 9, the two reference modes (Fig. 7a and b) are very roughly recovered by the PT-SGR trial. Furthermore, the “left” reference mode (Fig. 7a) appears to be better identified (compare Fig. 7a with Fig. 9c, d, f, and i). This finding is fairly positive given the relatively limited computational budget and the use of a small temperature ladder in regard to the problem dimensionality. Perhaps not surprisingly, the SGR performance is substantially worse. Here every independent chain is clearly trapped in one local optimum, which always has a larger data misfit than the reference RMSE of  $0.0030 \text{ kg/m}^{-3}$  (Fig. 10). Visual inspection of the final states of the 24 SGR trials also shows that almost none of the model realizations looks similar to any of the two reference modes (see Fig. 10 for three such examples). In average over the 24 trials, only 9% of the simulated pixels present a different facies between iterations 10,000 and 25,000. As a result, the mean ACF of SGR takes a value as high as 0.82 at lag 5000 (Fig. 11). With an ACF

value of 0.13, PT-SGR produces 6 times less autocorrelated samples at lag 5000 (Fig. 11).

#### 4. Discussion

Our results demonstrate that the standard SGR approach cannot cope efficiently with situations where the measurement data are collected at a relatively high spatial and/or temporal density. Standard SGR has however been shown to work for a data-poor situation, where the information content of the data does not constrain much the facies distribution and the posterior uncertainty is thus quite large [e.g., [22]].

Parallel tempering improves the SGR performance. Sampling of the complicated bimodal posterior distribution of case study 2 is hence much improved by parallel tempering, but not to the point of drawing samples from the correct stationary distribution within the allowed computing time and when using 24 temperatures (and thus parallel cores). Significantly increasing the number of temperatures, say by a factor 10, is expected to strongly enhance posterior sampling. Our future work will focus on two alternatives to simply increasing the available computing power. First, parallel tempering could be coupled with Wang–Landau (WL) sampling [1] for better exploration capabilities. The main principle of WL sampling

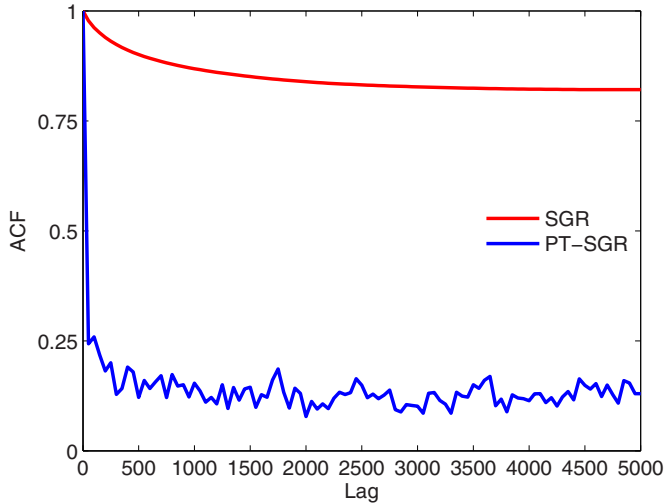


**Fig. 10.** Mean sampled model over MCMC iterations 10,000–25,000 (using a thinning factor of 50), and sampled models after 10,000 and 25,000 MCMC iterations for 3 out of the 24 independent SGR chains and case study 2. The red number in the lower left corner of each plot is the corresponding RMSE (m). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

is to adaptively bias the Metropolis acceptance probability in order to sample a flat histogram of posterior density values with pre-defined bins. The derived histogram thus contains an approximately equal amount of samples for each class of density value, and these samples can then be reused to approximate the posterior distribution (e.g., via importance sampling or by seeding a new MCMC run). The method might however not help with the observed problem of sampling slightly too large (log-)likelihoods, and thus a wrong stationary distribution. Second and perhaps more promising is the use of a more informal ensemble-based multiscale approach. The latter would consist in sequentially solving the inference problem from an upscaled coarser scale to the finer scale of interest, in the spirit of the work by Gardet et al. [7] for the multi-Gaussian case. The underlying idea is that the (upscaled) parameter space can be scanned quickly using a coarse resolution, thereby allowing for the subsequent finer scale searches to concentrate on the most productive parts of the prior distribution. Starting from different random points would then eventually provide an ensemble of solutions that (informally) approximate the posterior target.

On a more practical level, the rationale for our DS settings deserves special attention. A neighborhood of 50 (case study 1) or 75 (case study 2) nodes may seem large [24] as the DS simula-

tion time increases with the number of neighboring nodes. Nevertheless, using such large values was necessary to minimize the occurrence of artifacts in the generated models, which is caused by repeatedly conditioning on a large amount of grid points throughout the MCMC sampling. Combined with a large fraction of conditioning data (say  $> 50\%$  of the image), a small neighborhood can sometimes result in model proposals that are somewhat degraded compared to the TI. It is indeed the restricted neighborhood size that gives freedom to the DS to produce structures that are different than those found in the TI. Also, all it takes for a slightly degraded model to appear in the Markov chain is for the Metropolis rule (Eq. 5) to accept it. In other words, even if a model proposal with some artifacts is only rarely proposed, this model can persist in the Markov chain if the associated simulated data fit the observations sufficiently well. Even with the employed neighborhood of 50 nodes for case study 1 (see Section 3.2), trials  $\phi_{S1} = 0.5$  and  $\phi_{S1} = 0.75$  of Section 3.3 nevertheless showed artifacts in the proposed models, typically manifested by overly broad channels (not shown). To a lesser extent, other artifacts such as isolated patches and broken channels also occurred for trials  $\phi_{S2} = 0.9$  and  $\phi_{S2} = 0.95$  (not shown). The used distance threshold of 0.01 for case study 2 (see Section 3.8) also incurs a larger computational cost than that of using the more common value of 0.05 [e.g., [22]].



**Fig. 11.** Mean autocorrelation function (ACF) of the 7575 conductivity grid values derived from PT-SGR (blue line) or SGR (red line) for lags 0–5000 and case study 1. The lag- $k$  autocorrelation is defined as the correlation between draws  $k$  lags apart. Listed statistics are computed for the last 15,000 iterations of the unit temperature chain of PT-SGR or the 24 independent SGR chains, using a thinning factor of 50 thereby leading to a set of 300 sampled models for each chain. For SGR, the average of the 24 chains is presented. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article).

The value of 0.01 turned out to be required to (almost) systematically honor the point measurement data (see Section 3.7).

Finally, it would be interesting to investigate the performance of parallel tempering when used in conjunction with the patch-based geostatistical resimulation algorithm by Zahner et al. [42] which uses graph cuts. This method has been shown to be 40 times faster than DS for generating a 2D model, with a resulting posterior distribution that is (at least) of similar quality as that obtained by using DS.

## 5. Conclusion

This study is concerned with the application of sequential geostatistical resampling (SGR) to high-dimensional categorical field inference problems that present realistically complex likelihood functions. We highlight the limitations of the classical SGR approach and propose a parallel tempering implementation that, for a similar multi-core computing budget, provides much improved results with respect to both convergence towards the appropriate data misfit and sampling diversity. Two synthetic case studies are considered: a steady-state flow and a transport inverse problem, involving from 7501 to 10,000 unknowns. For the transport problem, the corresponding likelihood function is made bimodal with two well separated modes. In both case studies, every SGR MCMC chain gets trapped in a local optima while parallel tempering within sequential geostatistical resampling (PT-SGR) does not. The advantage of PT-SGR becomes more apparent for the bimodal inverse transport problem, for which PT-SGR is found to converge towards the reference data misfit much faster than SGR and to indicate the existence of two posterior modes. In contrast, for the same computational resources SGR appears to be barely able to appropriately fit the data and does almost not produce any single solution that looks visually similar to one of the two reference modes. Although PT-SGR outperforms SGR, our results also demonstrate that using a reasonably small number of temperatures (and thus parallel cores) in the range 16–24 may not allow sampling of the posterior distribution by PT-SGR within an affordable computational time. As an alternative to significantly increasing the num-

## Algorithm 1 Parallel tempering sequential geostatistical resampling.

```

1: procedure PT-SGR( $T, m, n, p_{SW}, meth_{SW}, P$ )  $\triangleright T$  is
   the temperature ladder of size  $n$  (with  $n$  even),  $m$  is the num-
   ber of MCMC iterations,  $p_{SW}$  is the probability of performing
   a swap update with temperature selection procedure  $meth_{SW}$ ,
   and  $P$  encapsulates the SGR algorithmic settings (e.g.,  $\phi$ , ...).
2:   for  $i = 1, \dots, m$  do  $\triangleright$  Loop over MCMC iterations
3:     for  $j = 1, \dots, n$  do in parallel  $\triangleright$  Loop over temperature
       ladder
4:        $p(\theta, T_j | \mathbf{d})_i \leftarrow \text{SGR}_{\text{MCMC}}(T_j, P)$   $\triangleright$  classical MCMC move
       with SGR. The  $j = 1, \dots, n$  updates are done in parallel.
5:     end for
6:     if  $p_{SW} > U(0, 1)$  then  $\triangleright$  Perform
       a swap update with probability  $p_{SW}$ , with  $U(0, 1)$  indicating an
       uniform random draw between 0 and 1.
7:        $\mathbf{r} = \text{SelectPairs}(n, meth_{SW}, i)$   $\triangleright$  Select pairs of
       temperatures
8:       for  $j = 1, \dots, n/2$  do
9:          $\mathbf{v} = \mathbf{r}(2 \times (j - 1) + 1)$     $\mathbf{w} = \mathbf{r}(2 \times j)$   $\triangleright$ 
       Propose swapping of selected pairs of chains, one possible ex-
       change swap per temperature
10:         $\alpha(\mathbf{v}, \mathbf{w}) \leftarrow 1 \wedge \left[ \frac{L(\theta_{\mathbf{w}} | \mathbf{d})}{L(\theta_{\mathbf{v}} | \mathbf{d})} \right]^{(1/T_{\mathbf{v}} - 1/T_{\mathbf{w}})}$ 
11:        if  $\alpha < U(0, 1)$  then  $\triangleright$  Swap chain temperatures
12:           $T_{\mathbf{v}} \leftarrow T_{\mathbf{w}}, \quad T_{\mathbf{w}} \leftarrow T_{\mathbf{v}}$ 
13:        end if
14:      end for
15:    end if
16:  end for
17: end procedure

```

## Algorithm 2 Selection of temperature indices at swapping time.

```

1: procedure  $\mathbf{R} = \text{SELECTPAIRS}(n, meth_{SW}, i)$ 
2:   if  $meth_{SW} = \text{"random"}$  then  $\triangleright$  Select temperatures randomly
3:      $\mathbf{r} = \text{permute}(n)$   $\triangleright$  Create a random permutation of the
       temperature indices
4:   elseif  $meth_{SW} = \text{"adjacent"}$   $\triangleright$  Consider all adjacent pairs
5:     if  $\text{mod}(i) = 1$  then  $\triangleright$  MCMC iteration number is odd
6:        $\mathbf{r} = [1, \dots, n]$ 
7:     else  $\triangleright$  MCMC iteration number is even
8:        $\mathbf{r} = [2, \dots, n - 1, 1, n]$ 
9:     end if
10:  end if
11: end procedure

```

ber of temperatures and thus computational needs, coupling PT-SGR with Wang–Landau sampling and (2) reframing SGR within an ensemble-based multiscale optimization framework are two potentially useful approaches that will be investigated in future work. More generally, PT could also prove useful when used in conjunction with dimensionality reduction approaches.

## Acknowledgments

A MATLAB code of the proposed PT-SGR approach is available from the first author. The patented DeeSee (DS) multiple-point statistics code is available for academic use upon request to one of its developers (Grégoire Mariethoz, Philippe Renard, Julien Straubhaar). Also, the synthetic measurements and forward modeling setup of our two case studies can be downloaded from <http://www.minds.ch/gm/downloads.htm>. We would like to thank the



associated editor and three anonymous referees for their positive feedback and useful comments.

## References

- [1] Bornn L, Jacob PE, Del Moral P, Doucet A. An adaptive interacting wang-landau algorithm for automatic density exploration. *J Comput Graph Stat* 2013;22(3):749–73. <http://dx.doi.org/10.1080/10618600.2012.723569>.
- [2] Carter JN, White DA. History matching on the imperial college fault model using parallel tempering. *Comput Geosci* 2013;17:43–65. <http://dx.doi.org/10.1007/s10596-012-9313-3>.
- [3] Earl DJ, Deem MW. Parallel tempering: theory, applications, and new perspectives. *Phys Chem Chem Phys* 2005;7:3910–16.
- [4] Falcioni M, Deem MW. A biased monte carlo scheme for zeolite structure solution. *J Chem Phys* 1999;110:1754–66.
- [5] Fu J, Gómez-Hernández JJ. A blocking markov chain monte carlo method for inverse stochastic hydrogeological modeling. *Math Geosci* 2009;41:105–28. <http://dx.doi.org/10.1007/s11004-008-9206-0>.
- [6] Galli A, Gao H. Rate of convergence of the gibbs sampler in the gaussian case. *Math Geol* 2001;33(6):653–77.
- [7] Gardet C, Le Ravalec M, Gloaguen E. Multiscale parameterization of petrophysical properties for efficient history-matching. *Math Geosci* 2014;46(3):315–36. <http://dx.doi.org/10.1007/s11004-013-9480-3>.
- [8] Gelman AG, Rubin DN. Inference from iterative simulation using multiple sequences. *Stat Sci* 1992;7:457–72.
- [9] Geman S, Geman D. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Trans Pattern Anal Mach Intell* 1984;6(6):721–41.
- [10] Geyer CJ. Markov Chain Monte Carlo maximum likelihood. *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*. New York: American Statistical Association; 1991. p. 156–63.
- [11] Geyer CJ. Chapter 11: Importance sampling, simulated tempering and umbrella sampling. In: Brooks S, Gelman A, Jones GL, Meng XL, editors. *Handbook of Markov Chain Monte Carlo*. Boca Raton: Chapman & Hall/CRC; 2011. p. 295–306.
- [12] Hansen TM, Mosegaard K, Cordua KC. Using geostatistics to describe complex a priori information for inverse problems. In: Ortiz JM, Emery X, editors. *VIII International Geostatistics Congress*. I. Santiago: Mining Engineering Department, University of Chile; 2008. p. 329–38.
- [13] Hansen TM, Cordua KC, Mosegaard K. Inverse problems with non-trivial priors: efficient solution through sequential gibbs sampling. *Comput Geosci* 2012;16:593–611. <http://dx.doi.org/10.1007/s10596-011-9271-1>.
- [14] Hansen TM, Cordua KC, Looms MC, Mosegaard K. SIPPI: a Matlab toolbox for sampling the solution to inverse problems with complex prior information: Part 2 - Application to crosshole GPR tomography. *Comput Geosci* 2013;52:481–92. <http://dx.doi.org/10.1016/j.cageo.2012.10.001>.
- [15] Hu LY. Gradual deformation and iterative calibration of gaussian related stochastic models. *Math Geol* 2000;32(1):87–108.
- [16] Khaninezhad MM, Jafarpour B, Li L. Sparse geologic dictionaries for subsurface flow model calibration: part i. inversion formulation. *Adv Water Resour* 2012;39(0):106–21. <http://dx.doi.org/10.1016/j.advwatres.2011.09.002>.
- [17] Künze R, Lunati I. An adaptive multiscale method for density-driven instabilities. *J Comput Phys* 2012;231:5557–70. <http://dx.doi.org/doi:10.1016/j.jcp.2012.02.025>.
- [18] Laloy E, Vrugt JA. High-dimensional posterior exploration of hydrologic models using multiple-try DREAM<sub>(ZS)</sub> and high-performance computing. *Water Resour Res* 2012;48(1). <http://dx.doi.org/10.1029/2011WR010608>.
- [19] Laloy E, Linde N, Jacques D, Vrugt JA. Probabilistic inference of multi-gaussian fields from indirect hydrological data using circulant embedding and dimensionality reduction. *Water Resour Res* 2015;51. <http://dx.doi.org/10.1002/2014WR016395>.
- [20] Lantuéjoul C. *Geostatistical simulation: models and algorithms*. Springer; 2002.
- [21] Lochbühler T, Vrugt JA, Sadegh M, Linde N. Summary statistics from training images as prior information in probabilistic inversion. *Geophys J Int* 2015;201:157–71. <http://dx.doi.org/10.1093/gji/ggv008>.
- [22] Mariethoz G, Renard P, Caers J. Bayesian inverse problem and optimization with iterative spatial resampling. *Water Resour Res* 2010a;46:W11530. <http://dx.doi.org/10.1029/2010WR009274>.
- [23] Mariethoz G, Renard P, Straubhaar J. The direct sampling method to perform multiple-point geostatistical simulations. *Water Resour Res* 2010b;46:W11536. <http://dx.doi.org/10.1029/2008WR007621>.
- [24] Meerschman E, Pirot G, Mariethoz G, Straubhaar J, Van Meirvenne M. A practical guide to performing multiple-point statistical simulations with the direct sampling algorithm. *Comput Geosci* 2013;52:307–24. <http://dx.doi.org/doi:10.1016/j.cageo.2012.09.019>.
- [25] Metropolis N, Rosenbluth AW, Rosenbluth MN, Teller AH, Teller E. Equation of state calculations by fast computing machines. *J Chem Phys* 1953;21:1087–92.
- [26] Mohamed L, Calderhead B, Filippone M, Christie M, Girolami M. Population MCMC methods for history matching and uncertainty quantification. *Comput Geosci* 2012;16(2):423–36. <http://dx.doi.org/10.1007/s10596-011-9232-8>.
- [27] Mosegaard K, Tarantola A. Monte carlo sampling of solutions to inverse problems. *J Geophys Res* 1995;100(B7):12431–47.
- [28] Opps SB, Schofield J. Extended state-space monte carlo methods. *Phys Rev E* 2001;60:056701. <http://dx.doi.org/10.1103/PhysRevE.63.056701>.
- [29] Predescu C, Predescu M, Ciobanu CV. On the efficiency of exchange in parallel tempering monte carlo simulations. *J Phys Chem B* 2005;109:4189–96. <http://dx.doi.org/10.1021/jp045073+>.
- [30] Rathore N, Chopra M, de Pablo JJ. Optimal allocation of replicas in parallel tempering simulations. *J Chem Phys* 2005;122:024111. <http://dx.doi.org/10.1063/1.1831273>.
- [31] Robert CP, Casella G. *Monte Carlo statistical methods*. 2nd ed. Springer; 2004.
- [32] Roberts GO, Gelman A, Gilks WR. Weak convergence and optimal scaling of random walk metropolis algorithms. *Ann Appl Probab* 1997;7:110–20.
- [33] Roberts GO, Rosenthal JS. Optimal scaling for various Metropolis-Hastings algorithms. *Stat Sci* 2001;16:351–67.
- [34] Romary T. Integrating production data under uncertainty by parallel interacting markov chains on a reduced dimensional space. *Comput Geosci* 2009;13:103–22. <http://dx.doi.org/10.1007/s10596-008-9108-8>.
- [35] Ruggeri P, Irving J, Holliger K. Systematic evaluation of sequential geostatistical resampling within MCMC for posterior sampling of near-surface geophysical inverse problems. *Geophys J Int* 2015;202:961–75. <http://dx.doi.org/10.1093/gji/ggv196>.
- [36] Sambridge M. A parallel tempering algorithm for probabilistic sampling and multimodal optimization. *Geophys J Int* 2014;196:357–74. <http://dx.doi.org/10.1093/gji/ggt342>.
- [37] Schrek A, Fort G, Moulines E. Adaptive equi-energy sampler: convergence and illustration. *ACM Trans Model Comput Simul* 2013;23(1):5:1–5:27. <http://doi.acm.org/10.1145/2414416.2414421>.
- [38] Strebelle S. Conditional simulation of complex geological structures using multiple point statistics. *Math Geol* 2002;34(1):1–22.
- [39] ter Braak CJ, Vrugt JA. Differential evolution markov chain with snooker updater and fewer chains. *Stat Comput* 2008;18(4):435–46. <http://dx.doi.org/10.1007/s11222-008-9104-9>.
- [40] Vo HX, Durlafsky LJ. A new differentiable parameterization based on principal component analysis for the low-dimensional representation of complex geological models. *Math Geosci* 2014;46:775–813. <http://dx.doi.org/10.1007/s11004-014-9541-2>.
- [41] Vrugt JA, Ter Braak C, Diks C, Robinson BA, Hyman JM, Higdon D. Accelerating markov chain monte carlo simulation by differential evolution with self-adaptive randomized subspace sampling. *Int J Nonlinear Sci Numer Simul* 2009;10(3):273–90. <http://dx.doi.org/doi:10.1515/ijnsns.2009.10.3.273>.
- [42] Zahner T, Lochbühler T, Mariethoz G, Linde N. Image synthesis with graph cuts: a fast model proposal mechanism in probabilistic inversion. *Geophys J Int* 2016;204(2):1179–90. <http://dx.doi.org/10.1093/gji/ggv517>.