

<sup>1</sup> **High-dimensional posterior exploration of hydrologic  
models using multiple-try DREAM<sub>(ZS)</sub> and high  
performance computing**

Eric Laloy,<sup>1</sup> and Jasper A. Vrugt<sup>12</sup>

---

E. Laloy, Department of Civil and Environmental Engineering, University of California, Irvine, California, USA. (elaloy@uci.edu).

J. A. Vrugt, Department of Civil and Environmental Engineering, University of California, Irvine, California, USA.

<sup>1</sup>Department of Civil and Environmental  
Engineering, University of California,  
Irvine, California, USA.

<sup>2</sup>Institute for Biodiversity and Ecosystems  
Dynamics, University of Amsterdam,  
Amsterdam, The Netherlands.

**Abstract.** Spatially distributed hydrologic models are increasingly being used to study and predict soil moisture flow, groundwater recharge, surface runoff, and river discharge. The usefulness and applicability of such complex models is increasingly held back by the potentially many hundreds (thousands) of parameters that require calibration against some historical record of data. The current generation of search and optimization algorithms is typically not powerful enough to deal with a very large number of variables, and summarize parameter and predictive uncertainty. In a previous paper [*Vrugt et al.*, 2008], we have presented a general-purpose Markov Chain Monte Carlo (MCMC) algorithm for Bayesian inference of the posterior probability density function of hydrologic model parameters. This method, entitled DiffeRelent Evolution Adaptive Metropolis (DREAM), runs multiple different Markov chains in parallel and uses a discrete proposal distribution to evolve the sampler to the posterior distribution. The DREAM approach maintains detailed balance and shows excellent performance on complex, multi-modal, search problems. Here, we present our latest algorithmic developments, and introduce MT-DREAM<sub>(ZS)</sub>, which combines the strengths of multi-try sampling, snooker updating, and sampling from an archive of past states. This new code is especially designed to solve high-dimensional search problems, and receives particularly spectacular performance improvement over other adaptive MCMC approaches when using distributed computing. Four different case studies with increasing dimensionality up to 241 parameters are used to illustrate the advantages of MT-DREAM<sub>(ZS)</sub>.

## 1. Introduction and scope

27 The ever increasing pace of computational power, along with significant advances in  
 28 measurement technologies, and interests in real-time forecasting has stimulated the devel-  
 29 opment of increasingly complex spatially distributed hydrologic models. The usefulness  
 30 and applicability of these models depends strongly on the values of the model parameters.  
 31 Unfortunately, the estimation of the correct values of these parameters have not proved to  
 32 be simple. To cope with the issues of heterogeneities, scale effects, and process complex-  
 33 ity, many models use effective parameters to aggregate complex interactions driven by a  
 34 number of highly interrelated energy and mass transport processes. The consequence of  
 35 this process aggregation is that the parameters cannot be directly measured in the field  
 36 but can only be meaningfully inferred by adjusting them so that the behavior of the model  
 37 approximates, as closely and consistently as possible, the observed system behavior over  
 38 some historical period of time. Sparse data and regionalization relationships may be used  
 39 to constrain the model by reducing the number of parameters, but the resulting inverse  
 40 problem nevertheless involves iterative improvement through successive executions of the  
 41 model, a situation that places a premium on calibration methods that can efficiently  
 42 summarize parameter and predictive uncertainty.

The past two decades have seen an increasing interest in Markov Chain Monte Carlo  
 methods for calibration of hydrologic models, and treatment of parameter, model struc-  
 tural, forcing data, and calibration data uncertainty [see, e.g., *Vrugt et al.*, 2003, 2008;  
*Keating et al.*, 2010; *Schoups and Vrugt*, 2010, and many others]. The basis of this method  
 is a Markov chain that generates a random walk through the search space and iteratively

finds solutions with stable frequencies stemming from a fixed probability distribution. To visit configurations with a stable frequency, an MCMC algorithm generates trial moves from the current position of the Markov chain at time  $t - 1$ ,  $\mathbf{x}_{t-1}$ , to a new state  $\mathbf{z}$ . The earliest MCMC approach is perhaps the well known random walk Metropolis (RWM) algorithm. Assuming that a random walk has already sampled the points  $\{\mathbf{x}_0, \dots, \mathbf{x}_{t-1}\}$ , this algorithm proceeds in the following three steps. First, a candidate point  $\mathbf{z}$  is sampled from a proposal distribution  $q(\cdot)$  that is symmetric,  $q(\mathbf{x}_{t-1}, \mathbf{z}) = q(\mathbf{z}, \mathbf{x}_{t-1})$  and may depend on the present location,  $\mathbf{x}_{t-1}$ . Next, the candidate point is either accepted or rejected using the Metropolis acceptance probability:

$$\alpha(\mathbf{x}_{t-1}, \mathbf{z}) = \begin{cases} \min\left[\frac{\pi(\mathbf{z})}{\pi(\mathbf{x}_{t-1})}, 1\right] & \text{if } \pi(\mathbf{x}_{t-1}) > 0 \\ 1 & \text{if } \pi(\mathbf{x}_{t-1}) = 0 \end{cases} \quad (1)$$

<sup>43</sup> where  $\pi(\cdot)$  denotes the density of the target distribution. Finally, if the proposal is  
<sup>44</sup> accepted, the chain moves to  $\mathbf{z}$ ; otherwise the chain remains at its current location  $\mathbf{x}_{t-1}$ .

The original RWM scheme was constructed to maintain detailed balance with respect to  $\pi(\cdot)$  at each step in the chain:

$$\pi(\mathbf{x}_{t-1})p(\mathbf{x}_{t-1} \rightarrow \mathbf{z}) = \pi(\mathbf{z})p(\mathbf{z} \rightarrow \mathbf{x}_{t-1}) \quad (2)$$

<sup>45</sup> where  $\pi(\mathbf{x}_{t-1})$  ( $\pi(\mathbf{z})$ ) denotes the probability of finding the system in state  $\mathbf{x}_{t-1}$  ( $\mathbf{z}$ ), and  
<sup>46</sup>  $p(\mathbf{x}_{t-1} \rightarrow \mathbf{z})$  ( $p(\mathbf{z} \rightarrow \mathbf{x}_{t-1})$ ) denotes the conditional probability of performing a trial move  
<sup>47</sup> from  $\mathbf{x}_{t-1}$  to  $\mathbf{z}$  ( $\mathbf{z}$  to  $\mathbf{x}_{t-1}$ ). The result is a Markov chain which, under certain regu-  
<sup>48</sup> larity conditions, has a unique stationary distribution with pdf  $\pi(\cdot)$ . In practice, this  
<sup>49</sup> means that if one looks at the values of  $\mathbf{x}$  generated by the RWM that are sufficiently  
<sup>50</sup> far from the starting value, the successively generated parameter combinations will be  
<sup>51</sup> distributed with stable frequencies stemming from the underlying posterior pdf of  $\mathbf{x}$ ,  $\pi(\cdot)$ .

<sup>52</sup> Hastings [1970] extended Eq. (2) to include non-symmetrical proposal distributions, i.e.  
<sup>53</sup>  $q(\mathbf{x}_{t-1}, \mathbf{z}) \neq q(\mathbf{z}, \mathbf{x}_{t-1})$ , in which a proposal jump to  $\mathbf{z}$  and the reverse jump do not have  
<sup>54</sup> equal probability. This extension is called the Metropolis-Hastings algorithm (MH), and  
<sup>55</sup> has become the basic building block of many existing MCMC sampling schemes.

<sup>56</sup> Existing theory and experiments prove convergence of well-constructed MCMC schemes  
<sup>57</sup> to the appropriate limiting distribution under a variety of different conditions. Yet, in  
<sup>58</sup> practice the convergence rate is often disturbingly slow. This inefficiency is typically  
<sup>59</sup> caused by an inappropriate selection of the orientation and scale of the proposal distri-  
<sup>60</sup> bution,  $q(\mathbf{x}_{t-1}, \cdot)$ , used to generate transitions in the Markov chain. When the proposal  
<sup>61</sup> distribution is too wide, very few candidate points will be accepted, and the chain will not  
<sup>62</sup> mix properly and converge rather slowly to the posterior target distribution. On the other  
<sup>63</sup> hand, if the proposal distribution is too narrow, the chain will remain in close vicinity  
<sup>64</sup> of its current location, and it will require a very large number of iterations before the  
<sup>65</sup> entire posterior distribution has been sampled. The selection of the proposal distribution  
<sup>66</sup> is therefore crucial in determining the efficiency and practical applicability of MCMC  
<sup>67</sup> simulation.

<sup>68</sup> In the past decade, a variety of different approaches have been proposed to increase  
<sup>69</sup> the efficiency of MCMC simulation and enhance the original RWM and MH algorithms.  
<sup>70</sup> These approaches can be grouped into single, and multiple chain methods. Single chain  
<sup>71</sup> methods work with a single trajectory, and continuously adapt the covariance,  $\Sigma$  of a  
<sup>72</sup> Gaussian proposal distribution,  $q_t(\mathbf{x}_{t-1}, \cdot) = N_d(\mathbf{x}_{t-1}, s_d \Sigma)$  using the information con-  
<sup>73</sup> tained in the sample path of the chain,  $\Sigma = \text{Cov}(\mathbf{x}_0, \dots, \mathbf{x}_{t-1}) + \varepsilon \mathbf{I}_d$ . The variable  $s_d$   
<sup>74</sup> represents a scaling factor (scalar) that depends only on the dimensionality  $d$  of the prob-

lem,  $\mathbf{x}$  is a  $d$ -dimensional vector,  $\mathbf{I}_d$  signifies the  $d$ -dimensional identity matrix, and  $\varepsilon$  is  
 a small scalar that slightly inflates the actual covariance,  $\Sigma$  so that the entire parameter  
 space can theoretically be sampled. As a basic choice, the scaling factor is chosen to  
 be  $s_d = 2.4^2/d$  which is optimal for Gaussian target and proposal distributions [Roberts  
 et al., 1997], and  $\varepsilon = 10^{-6}$ . Examples of self-adaptive single chain methods include the  
 Adaptive Metropolis (AM) [Haario et al., 2001] and Delayed Rejection Adaptive Metropo-  
 lis (DRAM) algorithms [Haario et al., 2006]. Component-wise updating of  $\mathbf{x}$  [Haario et  
 al., 2005] is possible to increase efficiency of AM for high-dimensional problems (large  $d$ ).  
 In addition, for the special case of hierarchical Bayesian inference of hydrologic models,  
 Kuczera et al. [2010] recently proposed to tune  $\Sigma$  using a limited-memory multi-block  
 pre-sampling step, prior to a classical single block Metropolis run.

Multiple chain methods use different trajectories running in parallel to explore the pos-  
 terior target distribution. The use of multiple chains has several desirable advantages,  
 particularly when dealing with complex posterior distributions involving long tails, cor-  
 related parameters, multi-modality, and numerous local optima [Gilks et al., 1994; Liu  
 et al., 2000; ter Braak, 2006; ter Braak and Vrugt, 2008; Vrugt et al., 2009; Radu et  
 al., 2009]. The use of multiple chains offers a robust protection against premature con-  
 vergence, and opens up the use of a wide array of statistical measures to test whether  
 convergence to a limiting distribution has been achieved [Gelman and Rubin, 1992]. One  
 popular multi-chain method that has found widespread application and use in hydrology  
 is the Shuffled Complex Evolution Metropolis algorithm [SCEM-UA, Vrugt et al., 2003].  
 Numerical experiments on a diverse set of mathematical test functions have shown that  
 SCEM-UA works well in practice. Yet, SCEM-UA does not generate a perfectly reversible

98 Markov chain. The explicit removal of outlier trajectories and covariance updating step  
 99 violate detailed balance. This poses questions on whether SCEM-UA generates an exact  
 100 sample of the posterior distribution. With some simple modifications, SCEM-UA could  
 101 be made an exact sampler, but this is beyond the scope of the current paper. We there-  
 102 fore consider the more recent Differential Evolution Markov chain (DE-MC) method of  
 103 ter Braak [2006]. This method is relatively easy to illustrate and understand, and can  
 104 be coded in just a few lines. DE-MC uses differential evolution as genetic algorithm for  
 105 population evolution with a Metropolis selection rule to decide whether candidate points  
 106 should replace their parents or not.

In DE-MC,  $N$  different Markov chains are run simultaneously in parallel. If the state of a single chain is given by a single  $d$ -dimensional vector  $\mathbf{x}$ , then at each generation the  $N$  chains in DE-MC define a population  $\mathbf{X}$ , which corresponds to an  $N \times d$  matrix, with each chain as a row. Jumps in each chain  $i = \{1, \dots, N\}$  are generated by taking a fixed multiple of the difference of two randomly chosen members (chains) of  $\mathbf{X}$  (without replacement) with indexes  $r_1$  and  $r_2$ :

$$\mathbf{z}^i = \mathbf{x}_{t-1}^i + \gamma(\mathbf{x}_{t-1}^{r_1} - \mathbf{x}_{t-1}^{r_2}) + \boldsymbol{\epsilon}, \quad r1 \neq r2 \neq i \quad (3)$$

107 where  $\gamma$  is a user-defined scalar, and  $\boldsymbol{\epsilon}$  is drawn from a symmetric  $d$ -dimensional dis-  
 108 tribution with a small variance compared to that of the posterior, but with unbounded  
 109 support. The difference vector in Eq. (3) contains the desired information about the  
 110 scale and orientation of the target distribution,  $\pi(\mathbf{x}|\cdot)$ . By accepting each jump with the  
 111 Metropolis ratio,  $\alpha(\mathbf{x}_{t-1}, \mathbf{z}) = \min [\pi(\mathbf{z}^i|\cdot)/\pi(\mathbf{x}_{t-1}^i|\cdot), 1]$ , a Markov chain is obtained, the  
 112 stationary or limiting distribution of which is the posterior distribution. The proof of this  
 113 is given in ter Braak and Vrugt [2008] and Vrugt et al. [2008, 2009]. Because the joint pdf

of the  $N$  chains factorizes to  $\pi(\mathbf{x}^1|\cdot) \times \dots \times \pi(\mathbf{x}^N|\cdot)$ , the states  $\mathbf{x}^1 \dots \mathbf{x}^N$  of the individual chains are independent at any generation after DE-MC has become independent of its initial value. After this burn-in period, the convergence of a DE-MC run can thus be monitored with the  $\hat{R}$ -statistic of *Gelman and Rubin* [1992]. From the guidelines of  $s_d$  in Random Walk Metropolis, the optimal choice of  $\gamma = 2.4/\sqrt{2d}$ . Every 10<sup>th</sup> generation,  $\gamma = 1.0$  to facilitate jumping between different modes [*ter Braak*, 2006].

DE-MC solves an important practical problem in random walk Metropolis, namely that of choosing an appropriate scale and orientation for the jumping distribution. Earlier approaches such as (parallel) adaptive direction sampling [*Gilks et al.*, 1994; *Roberts and Gilks*, 1994; *Gilks and Roberts*, 1996] solved the orientation problem but not the scale problem. Vrugt and coworkers [*Vrugt et al.*, 2008, 2009] showed that the efficiency of DE-MC can be enhanced, sometimes dramatically, using self-adaptive randomized subspace sampling and explicit consideration of aberrant trajectories. This method, entitled DiffeRential Evolution Adaptive Metropolis (DREAM), maintains detailed balance and ergodicity and has shown to exhibit excellent performance on a wide range of model calibration studies [e.g., *Vrugt et al.*, 2008; *Dekker et al.*, 2010; *Laloy et al.*, 2010a, b; *Scharnagl et al.*, 2010].

Unfortunately, standard DREAM (DE-MC) requires at least  $N = d/2$  to  $d$  ( $N = 2d$ ) chains to be run in parallel. Running many parallel chains is a potential source of inefficiency, as each individual chain requires burn-in to travel to the posterior distribution. The lower the number of chains required, the greater the practical applicability of DREAM for computationally demanding posterior exploration problems. One device that enables using a smaller  $N$  is to generate jumps in Eq. (3) from past states of the different chains.

*ter Braak and Vrugt* [2008] incorporated this idea into DE-MC and showed by numerical simulation and real-world examples that this method works well up to  $d = 100$  using only  $N = 3$  chains. These findings inspired *Vrugt et al.* [2011] to create DREAM<sub>(ZS)</sub> that capitalizes on the advantages of DREAM for posterior exploration but generates candidate points in each individual Markov chain by sampling from an archive of past states. This has several practical and theoretical advantages. Most importantly, only a few parallel chains ( $N = 3 - 5$ ) are required for posterior sampling. This reduces burn-in, particularly for problems involving many parameters (large  $d$ ), thereby increasing sampling efficiency. Indeed, initial studies to date presented in *Vrugt et al.* [2011] have shown that DREAM<sub>(ZS)</sub> requires fewer function evaluations than DREAM to converge to the appropriate limiting distribution. In DREAM<sub>(ZS)</sub>, the states of the chains are periodically stored in an archive using a simple thinning rule. The size of this matrix steadily increases during sampling, but the relative growth decreases linearly with generation number. This diminishing adaptation of the transition kernel ensures convergence of the individual chains to the posterior distribution [*Roberts and Rosenthal*, 2007]. To increase the diversity of the proposals, DREAM<sub>(ZS)</sub> additionally includes a snooker updater with adaptive step size. The snooker axis runs through the states of two different chains, and the orientation of this jump is different from the parallel direction update utilized in DREAM. The algorithmic implementation of the snooker update within the context of DE-MC is described in *ter Braak and Vrugt* [2008].

Despite significant enhancements in the efficiency of MCMC methods, it remains typically difficult to solve very high-dimensional posterior exploration problems involving hundreds or thousands of parameters. The performance of optimization and search meth-

160 ods typically deteriorates exponentially with increasing dimensionality of the parameter space. In applied mathematics this phenomenon is also referred to as the *Curse*  
 161 of dimensionality. This term was coined by Richard E. Bellman, within the context of  
 162 dynamic programming. In this paper, we present a general framework for efficient inversion  
 163 of highly-parameterized models. This method, entitled MT-DREAM<sub>(ZS)</sub>, combines  
 164 the strengths of differential evolution, subspace exploration, sampling from past states,  
 165 snooker updating, and multiple-try Metropolis sampling [Liu *et al.*, 2000] to efficiently  
 166 explore high-dimensional posterior distributions. This novel approach maintains detailed  
 167 balance and ergodicity and takes maximum advantage of distributed computing resources.  
 168 Four case studies with increasing complexity are used to demonstrate the advantages of  
 169 MT-DREAM<sub>(ZS)</sub> over current state-of-the-art optimization and MCMC algorithms, in-  
 170 cluding the Parameter ESTimation Toolbox [PEST, Doherty, 2009], the Shuffled com-  
 171 plexes with Principal component analysis [SP-UCI, Chu *et al.*, 2010], DREAM [Vrugt *et*  
 172 *al.*, 2009] and DREAM<sub>(ZS)</sub> [Vrugt *et al.*, 2011].  
 173

174 This paper is organized as follows. Section 2 presents the key concepts of MT-  
 175 DREAM<sub>(ZS)</sub> and discusses how to incorporate this new MCMC method on a high perfor-  
 176 mance computing platform. In section 3, we evaluate our algorithm against other state-  
 177 of-art MCMC methods, for two known mathematical benchmark distributions involving  
 178 multi-modality and high-dimensionality. This is followed by a real-world case study con-  
 179 sisting of the calibration of the Sacramento Soil Moisture Accounting model (SAC-SMA)  
 180 using daily discharge data from the Leaf River in Mississippi. This study involves only  
 181 13 parameters, but illustrates the severity of the hydrologic model calibration problem.  
 182 We conclude our testing of MT-DREAM<sub>(ZS)</sub> with a CPU-efficient 241-parameter ground-

<sup>183</sup> water model. This model calibration problem has been described in details in *Keating*  
<sup>184</sup> *et al.* [2010], and is used to illustrate the superior search capabilities of MT-DREAM<sub>(ZS)</sub>.  
<sup>185</sup> Finally, section 4 draws conclusions about the presented work and discusses yet to be  
<sup>186</sup> conceived methodological developments that will further increase the efficiency of MT-  
<sup>187</sup> DREAM<sub>(ZS)</sub>.

## 2. Theory and parallel implementation

<sup>188</sup> Our method merges the strengths of differential evolution, sampling from past states,  
<sup>189</sup> snooker updating, randomized subspace exploration, and multiple-try Metropolis sam-  
<sup>190</sup> pling [Liu *et al.*, 2000] for efficient high-dimensional posterior exploration. The resulting  
<sup>191</sup> new code, MT-DREAM<sub>(ZS)</sub>, is an extension of DREAM<sub>(ZS)</sub> [Vrugt *et al.*, 2011], and is es-  
<sup>192</sup> pecially designed for parallel implementation on a distributed computing cluster. We first  
<sup>193</sup> describe multiple-try Metropolis sampling (MTM), and then continue with an detailed  
<sup>194</sup> algorithmic description of MT-DREAM<sub>(ZS)</sub>.

### 2.1. Multiple-try Metropolis in MCMC sampling

<sup>195</sup> For reasons stated earlier, it is particularly important to have an appropriate selection of  
<sup>196</sup> the proposal distribution,  $q(\mathbf{x}_{t-1}, \cdot)$ , used to generate candidate points in each individual  
<sup>197</sup> chain. Local moves (small jumps) have a higher chance of being accepted but explore only  
<sup>198</sup> a small region. Large jumps, on the contrary, cover a larger part of the search space, yet  
<sup>199</sup> are typically rejected. Liu *et al.* [2000] have introduced a general approach that directly  
<sup>200</sup> confronts this tradeoff by creating multiple different candidate points simultaneously in-  
<sup>201</sup> volving both small and large jumps. This multiple-try Metropolis (MTM) approach has  
<sup>202</sup> several desirable advantages, one of them that the mixing of Markov chains is significantly

enhanced. In this work, we use the so-called MTM(II) variant, which was found to be the most robust for a range of different posterior exploration problems [Liu *et al.*, 2000]. This MTM(II) approach assumes a symmetric proposal distribution,  $q(\cdot|\cdot)$ , and can be described as follows:

1. Draw  $k$  trials  $\mathbf{z}_1, \dots, \mathbf{z}_k$  from  $q_t(\mathbf{x}_{t-1}, \cdot)$  where  $\mathbf{x}_{t-1}$  of size  $1 \times d$  denotes the current state of the chain.
2. Compute the posterior density,  $\pi(\mathbf{z}_j)$ , of each of the  $k$  proposal points,  $j = 1, \dots, k$ .
3. Randomly select one candidate point,  $\mathbf{z}_j$  of  $\mathbf{z}_1, \dots, \mathbf{z}_k$  with probability proportional to  $\pi(\mathbf{z}_j)$ .
4. Draw  $\mathbf{x}_1^*, \dots, \mathbf{x}_{k-1}^*$  reference points from  $q_t(\mathbf{z}, \cdot)$  and set  $\mathbf{x}_k^* = \mathbf{x}_{t-1}$ .
5. Accept  $\mathbf{z}$  with probability

$$\alpha(\mathbf{x}_{t-1}, \mathbf{z}) = \min \left\{ 1, \frac{\pi(\mathbf{z}_1) + \dots + \pi(\mathbf{z}_k)}{\pi(\mathbf{x}_1^*) + \dots + \pi(\mathbf{x}_k^*)} \right\} \quad (4)$$

This sampling scheme satisfies detailed balance, and therefore results in a reversible Markov chain with  $\pi(\mathbf{x})$  as its stationary distribution [Liu *et al.*, 2000]. Numerical studies in the same paper have shown that MTM(II) is considerably more efficient than a traditional MH sampler. This is particularly inspiring considering the large amount of wasted samples. For each transition in each of the chain,  $2k - 1$  samples are created of which only 1 is selected, and compared against the density of the current state of the chain. The information contained in the other  $2k - 2$  points is simply thrown away, and can therefore be considered wasted. In response to this, Frenkel [2004] has proposed a MCMC sampling strategy that recycles information from such rejected states. This further in-

<sup>222</sup> creases the efficiency of posterior sampling, but this approach has not found widespread  
<sup>223</sup> implementation and use.

<sup>224</sup> The use of multiple proposals in Eq. (4) however, places a heavier demand on compu-  
<sup>225</sup> tational resources, particularly when each candidate point is evaluated sequentially. For  
<sup>226</sup> example, lets assume we use  $k = 5$ . For each transition in the Markov chain, this choice  
<sup>227</sup> requires  $5 + 4 = 9$  different evaluations of  $\pi(\mathbf{x})$ . Hence, the  $k = 5$  proposal points need  
<sup>228</sup> to be evaluated, together with  $k - 1 = 4$  different points of the reference set. This is  
<sup>229</sup> computationally rather demanding. For the same computational budget, a single chain  
<sup>230</sup> method is able to create 9 different transitions in the Markov chain! Indeed, our numerical  
<sup>231</sup> tests do not corroborate the findings of *Liu et al.* [2000], but illustrate (not shown herein)  
<sup>232</sup> that multiple-try RWM with a (multi)normal proposal distribution (MTRWMN) gener-  
<sup>233</sup> ally requires more function evaluations than standard RWMN to converge to the target  
<sup>234</sup> distribution. This finding is consistent with the results of *Murray* [2007] who pointed  
<sup>235</sup> out a critical deficiency in *Liu et al.* [2000, p. 128, section 6.1] and also showed that if a  
<sup>236</sup> similar proposal distribution is used, RWMN outperforms MTRWMN.

<sup>237</sup> In both these studies, the candidate points of the proposal and reference set have been  
<sup>238</sup> evaluated sequentially. However, nothing prevents us from evaluating the  $k$  different pro-  
<sup>239</sup> posal trials, followed by the  $k - 1$  reference points simultaneously in parallel. This should  
<sup>240</sup> significantly enhance the efficiency of MTM(II). For example, if these  $k$  different points are  
<sup>241</sup> jointly evaluated then, in theory, parallel MTM(II) should be about  $(2k - 1) / 2$  more ef-  
<sup>242</sup> ficient than its sequential counterpart. Distributed computing thus significantly enhances  
<sup>243</sup> the efficiency of posterior exploration, particularly when dealing with computationally  
<sup>244</sup> demanding forward models [*Vrugt et al.*, 2006]. Yet, the efficiency of parallel MTM(II)

remains essentially dependent on the choice of the proposal distribution used to generate transitions in each of the individual Markov chains. The MTM(II) method uses a rather simplistic and fixed (multivariate normal) proposal distribution,  $q_t(\cdot|\cdot)$  to generate the points of the proposal and reference set. This is a potential source of inefficiency, in particular if the proposal distribution is a poor approximation of the target distribution [Vrugt *et al.*, 2003, 2008].

We hypothesize that significant efficiency improvements can be made if the proposal distribution of MTM(II) is adaptively updated en route to the posterior target distribution. Such updating significantly enhances acceptance rate, and the speed at which the posterior distribution is explored. It therefore seems logical to merge the strengths of MTM(II) and DREAM and create a single MCMC sampler that combines automatic proposal updating with multi-try sampling and parallel computing to further enhance the efficiency of posterior sampling, and provide a general-purpose algorithm that can efficiently solve difficult and high-dimensional search and optimization problems. Unfortunately, standard DREAM requires at least  $N = d/2$  to  $d$  chains to be run in parallel. This is a potential source of inefficiency, particularly for high-dimensional problems. For instance, for  $k = 5$  and  $d = 100$  parameters, we would need at least 250 - 500 different computational nodes to take full advantage of MTM(II), and accelerate the efficiency of posterior sampling. Most computer clusters will not readily have available such a large number of processors.

To minimize computational requirements, we capitalize on recent developments in MCMC simulation and combine DREAM<sub>(ZS)</sub> [Vrugt *et al.*, 2011] with MTM(II). This new code, entitled multi-try DREAM<sub>(ZS)</sub> or abbreviated MT-DREAM<sub>(ZS)</sub>, uses DREAM<sub>(ZS)</sub>

and MTM(II) as its main building block, but creates multiple proposals simultaneously in each of the  $N$  chains by sampling from an archive of past states. Previous studies have shown that DREAM<sub>(ZS)</sub> achieves excellent sampling efficiencies for  $d$  up to 50 - 100 using only  $N = 3$  different chains. Thus, if we create  $k = 5$  different proposals in each individual chain, then MT-DREAM<sub>(ZS)</sub> would require only 15 different nodes for optimal performance. Indeed, this is a much lower number of nodes than would be required with MT-DREAM. The MT-DREAM<sub>(ZS)</sub> code is especially designed to solve complex, high-dimensional inverse problems and summarize model and parameter uncertainty. The next section provides a detailed algorithmic description of MT-DREAM<sub>(ZS)</sub>, followed by four different case studies with increasing complexity.

## 2.2. DiffeRential Evolution Adaptive Metropolis with Multiple-try Sampling From an Archive of Past States → MT-DREAM<sub>(ZS)</sub>

We now describe our new code, entitled MT-DREAM<sub>(ZS)</sub>, which uses MTM(II) and DREAM<sub>(ZS)</sub> as main building blocks.

Let  $\mathbf{Z} = [x_j^i]$  ( $i = 1, \dots, M_0; j = 1, \dots, d$ ) be a  $M_0 \times d$  matrix, hereafter also referred to as archive, containing  $M_0$  draws from the prior distribution,  $p_d(\mathbf{x})$  of the  $d$  parameters. Similarly, let  $\mathbf{X}$  be a  $N \times d$ -matrix defining the  $N$  initialized starting positions,  $\mathbf{x}^i$ ,  $i = 1, \dots, N$  of the parallel chains by drawing samples from  $p_d(\mathbf{x})$ ;  $N \ll M_0$ . Lastly, let  $T$  be the number of population evolution steps and  $k$  be the number of parallel proposals. The initial population  $[\mathbf{X}_t; t = 0]$  is translated into a sample from the posterior distribution,  $\pi(\mathbf{x})$  using the following pseudo code:

```

287   1. Set  $M \leftarrow M_0$ 
288   FOR  $m \leftarrow 1, \dots, T$  DO (POPULATION EVOLUTION)

```

<sup>289</sup> FOR  $i \leftarrow 1, \dots, N$  DO (CHAIN EVOLUTION: A. PROPOSAL STEP)

<sup>290</sup>

(i) Generate  $l = 1, \dots, k$  candidate points,  $\mathbf{z}^{l,i}$  in chain  $i$ ,

$$\mathbf{z}^{l,i} = \mathbf{x}^i + (\mathbf{1}_d + \mathbf{e}_d)\gamma(\delta, d') \left[ \sum_{j=1}^{\delta} \mathbf{x}^{r_1(j)} - \sum_{n=1}^{\delta} \mathbf{x}^{r_2(n)} \right] + \mathbf{e}_d \quad (5)$$

<sup>291</sup> where  $\delta$  signifies the number of pairs used to generate the proposal,  $\mathbf{x}^{r_1(j)}$  and  $\mathbf{x}^{r_2(n)}$  are  
<sup>292</sup> rows from the archive  $\mathbf{Z}$ ;  $r_1(j), r_2(n) \in \{1, \dots, M\}$  and  $r_1(j) \neq r_2(n)$ . The values of  $\mathbf{e}_d$   
<sup>293</sup> and  $\mathbf{e}_d$  are drawn from  $U_d(-b, b)$  and  $N_d(0, b^*)$  with  $b$  and  $b^*$  small compared to the width  
<sup>294</sup> of the target distribution, respectively, and the value of the jump-size,  $\gamma$  depends on  $\delta$   
<sup>295</sup> and  $d'$ , the number of dimensions that will be updated jointly (see next step).

(ii) Replace each element ( $j = 1, \dots, d$ ) of the  $l = 1, \dots, k$  parallel proposals  $z_j^{l,i}$  with  
 $x_j^i$  using a binomial scheme with probability  $1 - CR$ ,

$$z_j^{l,i} = \begin{cases} x_j^i & \text{if } U \leq 1 - CR, \quad d' = d' - 1 \\ z_j^{l,i} & \text{otherwise} \end{cases} \quad j = 1, \dots, d \quad (6)$$

<sup>296</sup> where  $CR$  denotes the crossover probability, and  $U \in [0, 1]$  is a draw from a uniform  
<sup>297</sup> distribution.

<sup>298</sup> (iii) Compute  $\pi(\mathbf{z}^{l,i})$  for each of the  $l = 1, \dots, k$  proposals

<sup>299</sup> (iv) Select  $\mathbf{z}^i$  among the  $k$  proposals with probability  $\pi(\mathbf{z}^i)$

<sup>300</sup> END FOR (CHAIN EVOLUTION: A. PROPOSAL STEP)

<sup>301</sup> FOR  $i \leftarrow 1, \dots, N$  DO (CHAIN EVOLUTION: B. REFERENCE STEP)

<sup>302</sup>

(v) Generate  $l = 1, \dots, k - 1$  reference points,  $\mathbf{x}^{*,l,i}$  in chain  $i$  using Eqs. [5] and [6]

but now centered around  $\mathbf{z}^i$ ,

$$\mathbf{x}^{*,l,i} = \mathbf{z}^i + (\mathbf{1}_d + \mathbf{e}_d)\gamma(\delta, d') \left[ \sum_{j=1}^{\delta} \mathbf{x}^{r_1(j)} - \sum_{n=1}^{\delta} \mathbf{x}^{r_2(n)} \right] + \boldsymbol{\epsilon}_d \quad (7)$$

303 (vi) Compute  $\pi(\mathbf{x}^{*,l,i})$  for  $l = 1, \dots, k - 1$  and set  $\mathbf{x}^{*,k,i} = \mathbf{x}^i$  and  $\pi(\mathbf{x}^{*,k,i}) = \pi(\mathbf{x}^i)$

(vii) Accept  $\mathbf{z}^i$  with modified Metropolis acceptance probability:

$$\alpha(\mathbf{x}^{*,1,i}, \dots, \mathbf{x}^{*,k,i}; \mathbf{z}^{1,i}, \dots, \mathbf{z}^{k,i}) = \min \left\{ 1, \frac{\pi(\mathbf{z}^{1,i}) + \dots + \pi(\mathbf{z}^{k,i})}{\pi(\mathbf{x}^{*,1,i}) + \dots + \pi(\mathbf{x}^{*,k,i})} \right\} \quad (8)$$

304 (viii) If accepted, move the chain to the candidate point,  $\mathbf{x}^i = \mathbf{z}^i$ , otherwise remain  
305 at the old location,  $\mathbf{x}^i$ .

306

307 END FOR (CHAIN EVOLUTION: B. REFERENCE STEP)

308 END FOR (POPULATION EVOLUTION)

309

310 2. Append  $\mathbf{X}$  to  $\mathbf{Z}$  after each  $K$  steps, and then update  $M \leftarrow M + N$ .

311 3. Compute the *Gelman and Rubin* [1992] convergence diagnostic,  $\hat{R}_j$ , for each dimension  
312  $j = 1, \dots, d$  using the last 50% of the samples in each chain.

313 4. If  $\hat{R}_j \leq 1.2$  for all  $j$ , stop and go to step 5, otherwise go to POPULATION EVO-  
314 LUTION.

315 5. Summarize the posterior pdf using  $\mathbf{Z}$  after discarding the initial and burn-in samples.

316

317 The MT-DREAM<sub>(zs)</sub> algorithm is similar to DREAM, but uses multiple-try Metropolis  
318 sampling from an archive of past states to generate candidate points in each individual  
319 chain. This novel MCMC code has four main advantages. First, sampling from the past

circumvents the requirement of using  $N = d$  for posterior exploration. Especially for high-dimensional problems with large  $d$ , this has been shown to speed-up convergence to a limiting distribution [*ter Braak and Vrugt*, 2008; *Vrugt et al.*, 2011]. Second, the parallel multi-proposal implementation (see Figure 1) increases the efficiency of posterior exploration and accelerates the speed of convergence. Third, unlike DREAM, outlier chains do not require explicit consideration and removal. At any time during the simulation, transition from aberrant trajectories to the modal region remain possible by sampling their own immediate past state from  $\mathbf{Z}$  in combination with  $\gamma = 1$ . The chance of such jumps increases with increasing length of  $\mathbf{Z}$ . This is highly desirable. Even during burn-in, the  $N$  trajectories simulated with MT-DREAM<sub>(ZS)</sub> therefore maintain detailed balance at every single step in the chain [see *ter Braak and Vrugt*, 2008; *Vrugt et al.*, 2011, for proofs of detailed balance and ergodicity with sampling from past states]. Finally, as the proposal jumps in DREAM<sub>(ZS)</sub> are generated from an archive of past states, MT-DREAM<sub>(ZS)</sub> does not require sequential updating of the individual chains  $i = 1, \dots, N$  as implemented in DE-MC and DREAM to ensure detailed balance. This is of great advantage in a multiprocessor environment because each proposal point can then be simultaneously evaluated on a different node. The official proof of reversibility of DE-MC and DREAM demands the chains to be updated sequentially which impairs parallelization. Finally, MT-DREAM<sub>(ZS)</sub> contains a snooker update to increase the diversity of the candidate points, details of which can be found in *ter Braak and Vrugt* [2008], and *Vrugt et al.* [2011].

To speed up convergence to the target distribution, MT-DREAM<sub>(ZS)</sub> estimates a probability density function of different  $CR$  values during burn-in so that the average jumping distance is maximized. In practice, the probability  $p_m$  of  $n_{CR}$  different crossover val-

ues,  $CR = m/n_{CR} \mid m = 1, \dots, n_{CR}$ , is estimated by maximizing the squared distance,  
 $\Delta = \sum_{i=1}^N \sum_{j=1}^d (\bar{\mathbf{x}}_{j,t}^i - \bar{\mathbf{x}}_{j,t-1}^i)^2$  between the two subsequent samples  $\bar{\mathbf{x}}_t$  and  $\bar{\mathbf{x}}_{t-1}$  of the  
 $N$  different chains. The position of the chains is normalized with the standard deviation  
of each individual dimension (parameter) calculated from the current population,  $\mathbf{X}_t$ , so  
that all  $d$  dimensions contribute equally to  $\Delta$ . The algorithm results in an optimized  
probability for each individual  $CR$  value. This distribution is determined during burn-in  
and used to randomly select a  $CR$  value for each different proposal point, which in turn  
determine the effective number of dimensions or  $d'$  used to calculate  $\gamma(\cdot, \cdot)$  in Eqs. (5)  
and (7). A detailed description of this adaptation strategy appears in *Vrugt et al.* [2009]  
and so will not be repeated herein.

### 2.3. THEOREM

*Suppose  $\pi(\cdot)^N$  is a fixed target probability distribution, on a state space  $\mathbf{X}$ . MT-DREAM<sub>(ZS)</sub> constructs a Markov chain kernel  $P(\cdot)$  which has  $\pi(\cdot)^N$  as its stationary distribution such that:*

$$\|P_T(\mathbf{x}, \cdot) - \pi(\cdot)^N\| \rightarrow_{T \rightarrow \infty} 0 \quad (9)$$

for any  $\mathbf{x} \in \mathbf{X}$ .

**Proof:** We are left with giving a formal proof of detailed balance of MT-DREAM<sub>(ZS)</sub>.  
Details of this can be found in *Vrugt et al.* [2011] for DREAM<sub>(ZS)</sub> and *Liu et al.* [2000] for  
MTM(II), and we refer the reader to these two papers. In few words, we conjecture that  
DREAM<sub>(ZS)</sub> yields a Markov chain that is ergodic with unique stationary distribution with  
pdf  $\pi(\cdot)^N$  because of its diminishing adaptation [*Roberts and Rosenthal*, 2007]. Indeed,  
the matrix  $\mathbf{Z}$  grows during the course of the sampling process by an order  $N/M = K/t$ ,

360 which decreases in generation time  $t$ . Adaptive changes in the proposal distribution (and  
 361 thus in the transition kernel  $P(\cdot)$ ) therefore diminish to zero as the length of the thinned  
 362 past increases. Because both DREAM<sub>(ZS)</sub> and MTM(II) generate a reversible Markov  
 363 chain, so thus their combination. This concludes our proof.

#### 2.4. Selection of Algorithmic Variables in the MT-DREAM<sub>(ZS)</sub> algorithm

364 The MT-DREAM<sub>(ZS)</sub> algorithm contains several algorithmic variables that need to be  
 365 specified by the user before the method can be used for posterior inference. These variables  
 366 include  $N$ , the number of chains,  $k$  the number of parallel proposal trials in each individual  
 367 chain,  $M_0$ , the initial size of the archive  $\mathbf{Z}$ , the thinning rate  $K$  used to periodically record  
 368 samples in  $\mathbf{Z}$ , and the probability of performing a snooker update [see *ter Braak and*  
 369 *Vrugt, 2008; Vrugt et al., 2011*, for details],  $p_{SK}$ . Based on recommendations in *Vrugt*  
 370 *et al.* [2011], we set  $N = 3 - 5$ ,  $K = 10$ ,  $M_0 = 10d$  and  $p_{SK} = 0.1$ , and use default  
 371 values of  $\delta = 1$ ,  $b = 0.05$ ,  $b^* = 10^{-6}$ , and  $n_{CR} = 3$  from the original two DREAM  
 372 papers *Vrugt et al.* [2008, 2009]. From the guidelines of  $s_d$  in RWM, the optimal choice of  
 373  $\gamma(\delta, d') = 2.38/\sqrt{2\delta d'}$  in Eqs. (5) and (7). To help facilitate direct jumps between different  
 374 disconnected posterior modes, we temporarily switch to  $\gamma = 1$  at every 5<sup>th</sup> proposal point  
 375 [*Vrugt et al.*, 2008, 2009].

376 This leave us with choosing a value for  $k$ . A few preliminary tests for a range of different  
 377 search problems suggests that  $k = 5$  works well in practice. We therefore recommend this  
 378 value in future applications.

### 3. Cases studies

379 To illustrate the efficiency of MT-DREAM<sub>(ZS)</sub>, we conducted a wide range of numerical  
 380 experiments. These tests include two known mathematical target distributions, and two  
 381 real-world studies involving the calibration of the Sacramento Soil Moisture Accounting  
 382 (SAC-SMA) model [Burnash, 1995], and a 241-parameter groundwater model. These case  
 383 studies cover a diverse set of problem features, including high-dimensionality, nonlinear-  
 384 ity, non-convexity, multi-modality, and numerous local optima. In all our calculations  
 385 with MT-DREAM<sub>(ZS)</sub>, we use the default settings of the algorithmic variables specified  
 386 previously. We use  $N = 3$  parallel chains, but temporarily switch to  $N = 5$  for one  
 387 of the case studies involving significant multi-modality. To benchmark the results of  
 388 MT-DREAM<sub>(ZS)</sub>, we include comparison against the DREAM [Vrugt *et al.*, 2008, 2009],  
 389 DREAM<sub>(ZS)</sub> [Vrugt *et al.*, 2011], and RWMN algorithms using standard settings of the  
 390 algorithmic variables reported in the literature. Note that the RWMN sampler runs only a  
 391 single chain, and uses a multivariate normal proposal distribution,  $N_d(0, c\mathbf{I}_d)$  with  $c$  tuned  
 392 to get an acceptance rate of about 24%. This is typically considered optimal [Roberts *et*  
 393 *al.*, 1997].

394 The algorithmic developments presented in this paper have been inspired by the in-  
 395 creasing availability of distributed computing resources. Indeed, the potential of parallel  
 396 computing sheds a completely different light on what constitutes an efficient algorithm.  
 397 Widely celebrated search and optimization algorithms that have received a lot of attention  
 398 in the past decades, might no longer be most efficient in a multi-processor environment.  
 399 Their inherent sequential topology limits multi-tasking. An example of this includes  
 400 the single-chain AP [Haario *et al.*, 1999], AM [Haario *et al.*, 2001] and DRAM methods  
 401 [Haario *et al.*, 2005]. Their current sequential topology prevents effective use of distributed

computing resources. Multi-chain methods on the contrary are easier to parallelize but it is important to preserve detailed balance. This requirement dictates that the  $N$  different chains in DREAM are updated sequentially, and we therefore run this algorithm on a single processor. The DREAM<sub>(zs)</sub> and MT-DREAM<sub>(zs)</sub> algorithms, on the contrary are specifically designed for implementation on a distributed computing network. Sampling from the past ensures reversibility even if the  $N$  chains and / or multiple candidate points are evaluated jointly in parallel. We therefore execute DREAM<sub>(zs)</sub> and MT-DREAM<sub>(zs)</sub> in parallel using multiple different processors.

The differences in computer implementation of the DREAM, DREAM<sub>(zs)</sub> and MT-DREAM<sub>(zs)</sub> algorithms complicates a comparative efficiency analysis. For instance, for the computational costs of a single proposal evaluation in DREAM, the DREAM<sub>(zs)</sub> algorithm is able to execute  $N$  different candidate points simultaneously in parallel. Widely used measures such as the total number of function evaluations, hereafter referred to as  $FE$  is thus no longer sufficient to compare the efficiency of the different MCMC algorithms used herein. We therefore introduce a Computational Time Unit or  $CTU$ . Sequential MCMC samplers, such as RWMN and DREAM, use one  $CTU$  for each  $FE$ , thus essentially  $CTU = FE$ . Parallel samplers, on the contrary, evaluate multiple proposal points in parallel, thus automatically  $CTU < FE$ . In particular, for DREAM<sub>(zs)</sub>, it is not difficult to demonstrate that  $CTU = FE/N$ . This leaves us with MT-DREAM<sub>(zs)</sub>. Unfortunately, it is not immediately obvious what the mathematical relationship is between  $CTU$  and  $FE$  for this particular algorithm. The proposal and reference steps (Eqs. (5)) and (7)) both need a single  $CTU$  but involve a different number of function evaluations. The proposal step simultaneously evaluates  $N \times k$  points (i.e.  $FE$ ), whereas the reference

425 step executes  $N \times (k - 1)$  candidate points in parallel. A single  $CTU$  in MT-DREAM<sub>(ZS)</sub>  
 426 is thus equivalent to approximately  $N \times (k - \frac{1}{2}) FE$ , which after rearrangement gives  
 427  $CTU = FE/(N \times (k - \frac{1}{2}))$ . Also, because the proposal and reference step in MT-  
 428 DREAM<sub>(ZS)</sub> require 2  $CTU$ , MT-DREAM<sub>(ZS)</sub> is theoretically about 50% less efficient  
 429 than DREAM<sub>(ZS)</sub>! In practice, however the multi-try step will exhibit some important  
 430 advantages, as will be demonstrated later. To make the comparison between DREAM<sub>(ZS)</sub>  
 431 and MT-DREAM<sub>(ZS)</sub> as fair as possible, our numerical experiments presented herein also  
 432 include DREAM<sub>(ZS)</sub> with a number of chains similar to the number of parallel processors  
 433 using by MT-DREAM<sub>(ZS)</sub>, which is simply  $N \times k$ .

434 Note that these developments essentially ignore the time required for communication  
 435 between the master and slave nodes. In most practical applications involving CPU in-  
 436 tensive forward models, the time it requires to communicate between the master and the  
 437 slave nodes is negligibly small compared to the time it requires to execute the actual  
 438 simulation model, and compute the desired output. The most efficient MCMC method  
 439 requires the lowest number of  $CTU$  to generate samples from the posterior distribution.

440 We use three different diagnostic measures to check when convergence of each sampler to  
 441 a limiting distribution has been achieved. The first diagnostic is the  $\hat{R}$  statistic of *Gelman*  
 442 and *Rubin* [1992] which compares the between and within variance of the different chains.  
 443 Convergence is declared when  $\hat{R}_j \leq 1.2$  for all  $j = 1, \dots, d$ , and the corresponding  $CTU$   
 444 is denoted with  $CTU_{\hat{R}}$ .

445 The second and third convergence diagnostic are derived using the approach of *Raftery*  
 446 and *Lewis* [1992]. Suppose that we like to measure some posterior quantile, hereafter  
 447 referred to as *qnt*. If we define a tolerance *r* of *qnt* and a probability *s* of being within that

tolerance, the Raftery-Lewis diagnostic estimates the number of posterior samples,  $RL_{NT}$ , and the required burn-in length of the Markov chain,  $RL_{BURN}$ , necessary to satisfy the given tolerances. Yet, the Raftery-Lewis diagnostic will differ depending on what quantile is being chosen. We therefore follow *El Adlouni et al.* [2006] and compute  $RL_{BURN}$  and  $RL_{NT}$  for 9 different values of  $qnt \in 0.1, 0.2, \dots, 0.9$ . We do this for each  $j$ th dimension ( $j = 1, \dots, d$ ) of the posterior target, and retain the largest values of  $RL_{BURN}$  and  $RL_{NT}$ . We then report the number of *CTU* needed for the sampler to produce those required number of samples. We follow the statistical literature and set  $r = 0.9$ .

The final performance criteria considered herein measures the autocorrelation between the various samples of the Markov chains created with the different MCMC algorithms. We use the so-called Inefficiency Factor (*IF*) [*Chib et al.*, 2002]. In principle, the *IF* is similar to the inverse of the numerical efficiency measure of *Geweke* [1992] and can be computed from the last samples in the Markov chains [*Chib et al.*, 2002]. We compute this *IF* criterion for each dimension and report the largest value. To be able to compare the values of *IF* for the different MCMC methods, we estimate this criterion using a similar number of (final) posterior samples. This number of samples is simply taken to be 25% of the maximum number of *FE* of the sequential codes.

The four different performance criteria discussed so far primarily estimate the time required to reach convergence for each individual method, and generate high-quality (and thus uncorrelated) samples from the posterior distribution. These diagnostics essentially measure efficiency, without recourse to estimating the correctness of the sampled posterior distribution. For example, consider a case in which an algorithm has prematurely (quickly) converged to the wrong distribution. The convergence criteria would actually

convey a spectacular performance! We therefore need to augment these three different efficiency diagnostics with criteria that explicitly measure the distance to the actual target distribution. Of course, such effectiveness criteria can only be computed if the posterior distribution is actually known beforehand. We consider two of such synthetic distributions in this paper, and introduce an additional diagnostic measure,  $D$ , that measures the average normalized Euclidean distance to the true mean  $\mu_\pi$  and standard deviation  $\sigma_\pi$  of the posterior target distribution:

$$D = \sqrt{\frac{1}{2d} \sum_{i=1}^d \left[ \left( \frac{\mu_\pi - \hat{\mu}_\pi}{\sigma_\pi} \right)^2 + \left( \frac{\sigma_\pi - \hat{\sigma}_\pi}{\sigma_\pi} \right)^2 \right]} \quad (10)$$

where the superscript  $\hat{\cdot}$  denotes the posterior moments derived from the posterior draws generated with each sampler. We take a similar number of (final) draws for each different MCMC method, and this number is identical to 25% of the total number of *FE* allowed for the sequential samplers. A similar calculation is used for *IF*. For our mathematical test distributions, we report the values of  $D$  alongside with  $CTU_{\hat{R}}$ ,  $RL_{BURN}$ ,  $RL_{NT}$  and *IF*. In all our calculations reported herein, we do not consider the Delayed rejection Adaptive Metropolis algorithm [DRAM *Haario et al.*, 2006] because this adaptive sampler has shown to exhibit rather poor performance [*Vrugt et al.*, 2009].

### 3.1. A 200-dimensional multivariate normal distribution

To test the performance of our code in the presence of high-dimensionality, the first case study considers a 200-dimensional multivariate normal distribution, centered at the zero vector. The covariance matrix was set such that the variance of the  $j^{th}$  variable was equal to  $j$ , with pairwise correlations of 0.5. The initial population is drawn from  $\mathbf{X} \in [-5.0, 15.0]^d$  reflecting a lack of prior knowledge about the mean and variance of the

478 posterior. A maximum total of 1,000,000 *CTU* were allowed for the sequential RWMN and  
 479 DREAM methods. For the parallel DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub> codes a maximum  
 480 total of 400,000 *CTU* was deemed sufficient to explore the posterior target distribution.  
 481 Thus, RWMN and DREAM were allowed to use more than two times the amount of time  
 482 assigned to DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>.

483 Table 1 presents summary statistics of 25 subsequent trials for each of the four different  
 484 Metropolis samplers. The results presented in this table highlight several important find-  
 485 ings. First, MT-DREAM<sub>(ZS)</sub> provides the closest approximation of the actual target distri-  
 486 bution. Although, DREAM<sub>(ZS)</sub> with  $N = 15$  chains converges the fastest of all the different  
 487 algorithms, the resulting posterior samples not only exhibit considerably more autocorrela-  
 488 tion, but also are less consistent with the true posterior distribution. Second, multiple-try  
 489 sampling significantly enhances the mixing of the different Markov chains. Approximately  
 490 45% of the proposal points is being accepted with MT-DREAM<sub>(ZS)</sub>, whereas an accep-  
 491 tance rate of about 17 - 24% is found for the other MCMC samplers. This partly explains  
 492 the superior performance of MT-DREAM<sub>(ZS)</sub>. Note that acceptance rate of 45.2% ob-  
 493 tained with MT-DREAM<sub>(ZS)</sub> falls within the range of 30-50% recommended by *Liu et al.*  
 494 [2000] for standard MTM. Third, and as anticipated, notice the inferior results of RWMN.  
 495 The fixed proposal distribution (identity matrix), albeit scaled to receive an acceptance  
 496 rate of 24%, is a rather poor approximation of the actual target distribution, and hence  
 497 many *FE* or *CTU* are necessary with RWMN to approximate the posterior distribution.  
 498 Fourth, and as hypothesized in the introduction, the combination of parallel evalution  
 499 of the candidate points with sampling from the past in DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>

500 tremendously reduces the required burn-in. Lastly, MT-DREAM<sub>(ZS)</sub> generates posterior  
 501 samples with the smallest (auto)correlation among the different codes.

502 To provide insights into the sampled posterior distributions, please consider Figure 2  
 503 which presents histograms of dimension 1 and 200 of the multivariate normal distribu-  
 504 tion. The red lines represent the true marginal distributions of  $x_1$  and  $x_{200}$ . The RWM  
 505 samples receives particular poor performance, and the marginal distributions deviate con-  
 506 siderably from their true counterpart. Multi-chain methods (bottom three panels) receive  
 507 a noticeable better performance, but overall MT-DREAM<sub>(ZS)</sub> receives superior results.  
 508 The posterior distribution sampled with this method almost perfectly matches the actual  
 509 target distribution.

### 3.2. A 25-dimensional tri-modal distribution

510 The second case study involves a 25-dimensional tri-modal pdf with three disconnected  
 511 modes. This example builds on the bimodal distribution presented in *Vrugt et al.* [2009],  
 512 and is given by  $\pi(\mathbf{x}) = 3/6N_d(\mathbf{10}, \mathbf{I}_d) + 2/6N_d(\mathbf{5}, \mathbf{I}_d) + 1/6N_d(-\mathbf{5}, \mathbf{I}_d)$  where **10**, **5** and  
 513 **-5** are  $d$ -dimensional vectors. This distribution is notoriously difficult to approximate  
 514 with MCMC simulation, because the three individual modes are so far separated that  
 515 standard Metropolis samplers cannot jump from one mode to the other. This complicates  
 516 convergence. A maximum total of 2,000,000 *CTU* was deemed sufficient for the sequential  
 517 RWMN and DREAM algorithms to explore the target distribution. The DREAM<sub>(ZS)</sub> and  
 518 MT-DREAM<sub>(ZS)</sub> codes on the contrary were allowed to only use 400,000 *CTU*. This is still  
 519 sufficient to reach convergence to a limiting distribution. Thus, the parallel codes consume  
 520 only 1/5 of the total computational time that is assigned to RWMN and DREAM.

521 Table 2 summarizes the performance of RWMN, DREAM, DREAM<sub>(ZS)</sub> and MT-  
 522 DREAM<sub>(ZS)</sub>. Listed statistics denotes averages over 25 independent trials. The results  
 523 presented in this table are qualitatively very similar to those of the previous study, and  
 524 even more clearly highlight the excellent performance of MT-DREAM<sub>(ZS)</sub>. As expected,  
 525 RWMN exhibits a particular poor performance. Each different trial converges to a dif-  
 526 ferent posterior mode resulting in a rather poor approximation of the target distribution  
 527 and thus rather high value of the  $D$  statistic. A single chain is typically unable to cope  
 528 with this multi-modal search space, and provide an accurate characterization of the tar-  
 529 get distribution. As expected, significantly better results are obtained if multiple different  
 530 trajectories are run simultaneously. Not only, does the discrete proposal distribution of  
 531 Eq. (3) allow for immediate jumps between the different disconnected modes, the use of  
 532 a number of different chains also protects against premature convergence. Yet, standard  
 533 DREAM and DREAM<sub>(ZS)</sub> exhibit a rather poor acceptance rate. Less than 1 out of 10  
 534 candidate points is being accepted which causes a rather slow mixing of the individual  
 535 Markov chains. Much better results are obtained with MT-DREAM<sub>(ZS)</sub> when multiple  
 536 candidate points are jointly considered in each individual chain. This not only significantly  
 537 increases the acceptance rate to about 30% but also results in the closest approximation  
 538 of the target distribution.

539 Figure 3 presents the results of our analysis. The top 3 graphs present histograms of the  
 540 sampled  $x_1$  values using the (3a) DREAM, (3b) DREAM<sub>(ZS)</sub>, and (3c) MT-DREAM<sub>(ZS)</sub>  
 541 algorithm. The red line depicts the true marginal posterior pdf of the tri-modal test func-  
 542 tion. The results conclusively show that MT-DREAM<sub>(ZS)</sub> exhibits the best performance  
 543 and provides the closest approximation of the true target distribution. This finding is

544 consistent with the results reported in Table 2, and inspires confidence in the ability of  
 545 MT-DREAM<sub>(zs)</sub> to deal with multi-modal posterior distributions.

546 The bottom panel (3d) shows a trace plot of the sampled  $x_1$  values derived with MT-  
 547 DREAM<sub>(zs)</sub>. Each of the  $N$  Markov chains is coded with a different color. The different  
 548 parallel trajectories mix well, and jump back-and-forth between the different posterior  
 549 modes. The density of the points in each mode is consistent with the weight of each peak  
 550 in the actual target distribution. These findings highlight the ability of MT-DREAM<sub>(zs)</sub>  
 551 to efficiently explore multi-modal posterior distributions.

### 3.3. The rainfall - runoff transformation: SAC-SMA model

552 Our first real-world case study considers calibration of the Sacramento Soil Moisture  
 553 Accounting (SAC-SMA) model [Burnash, 1995]. The SAC-SMA model is a lumped con-  
 554 ceptual watershed model that describes the transformation from rainfall into basin runoff  
 555 using six different reservoirs (state variables). A unit hydrograph is commonly used to rout  
 556 channel inflow downstream and compute streamflow at the gauging point. This model  
 557 is extensively used by the National Weather Service for flood forecasting throughout the  
 558 United States, and has 13 user-specifiable (and 3 fixed) model parameters, which are listed  
 559 in Table 3. Inputs to the model include mean areal precipitation (MAP) and potential  
 560 evapotranspiration (PET) while the outputs are estimated evapotranspiration and chan-  
 561 nel inflow. Various studies have demonstrated that calibration of the SAC-SMA model is  
 562 very difficult due to the presence of numerous local optima in the parameter space with  
 563 both small and large domains of attraction, discontinuous first derivatives, and curving  
 564 multidimensional ridges [Duan *et al.*, 1992; Thiemann *et al.*, 2001; Vrugt *et al.*, 2006, 2009;

<sup>565</sup> *Chu et al.*, 2010]. Although this study only involves 13 different parameters, it poses an  
<sup>566</sup> interesting challenge for MCMC samplers, as will be shown later on.

<sup>567</sup> We estimate the posterior distribution of the SAC-SMA parameters using historical data  
<sup>568</sup> from the Leaf River watershed. This humid basin of approximately 1,950 km<sup>2</sup> is located  
<sup>569</sup> North of Collins, Mississippi, in the United States. We used 2 years of daily discharge data  
<sup>570</sup> from Jan. 1, 1953 to Dec. 31, 1954 to estimate the SAC-SMA parameters. In practice, it is  
<sup>571</sup> advisable to use a longer record of calibration data [*Yapo et al.*, 1996; *Vrugt et al.*, 2006],  
<sup>572</sup> but deliberately we use only a 2 year record to result in a very challenging parameter  
<sup>573</sup> estimation problem with multiple disconnected regions of attraction [*Vrugt et al.*, 2009].

<sup>574</sup> We assume a flat or uniform prior distribution of the SAC-SMA model parameters with  
<sup>575</sup> ranges specified in Table 3.

A simple Sum of Squared Error (SSE) loss function was used to compare the SAC-SMA predictions of daily streamflow with their respective observed discharges. With the assumption of a noninformative prior, the SSE results in the following posterior density function:

$$\pi(\mathbf{x}|\hat{\mathbf{y}}, \boldsymbol{\phi}) = \left[ \sum_{i=1}^{N_m} (y_i(\mathbf{x}, \boldsymbol{\phi}) - \hat{y}_i)^2 \right]^{-\frac{1}{2}N_m} \quad (11)$$

<sup>576</sup> where  $y_i(\mathbf{x}, \boldsymbol{\phi})$  ( $\hat{y}_i$ ) denotes the model predicted (observed) streamflow, respectively,  $N_m$   
<sup>577</sup> signifies the number of measurements, and  $\boldsymbol{\phi}$  represents the initial and forcing conditions.  
<sup>578</sup> A maximum total of 100,000 CTU was used to approximate the posterior distribution of  
<sup>579</sup> the 13 SAC-SMA model parameters.

<sup>580</sup> Our previous work [*Vrugt et al.*, 2009] has demonstrated the superiority of DREAM over  
<sup>581</sup> other adaptive MCMC samplers including the optimal RWMN sampler, DRAM [*Haario  
<sup>582</sup> et al.*, 2006] and DE-MC [*ter Braak*, 2006]. This work also illustrated a rather poor per-

583 formance of the Shuffled Complex Evolution (SCE-UA) global optimization algorithm of  
 584 *Duan et al.* [1992]. This method was originally developed in the early 1990s to solve highly  
 585 nonlinear, non-convex and non-continuous optimization problems, and because of its effi-  
 586 ciency and effectiveness has become the method of choice for watershed model calibration  
 587 problems. Yet, for this particular problem, SCE-UA was found to get stuck in a local  
 588 basin of attraction with RMSE values of about  $13.7 \text{ m}^3 \text{ s}^{-1}$ , whereas DREAM identified  
 589 a global minima around  $13.2 \text{ m}^3 \text{ s}^{-1}$ . This difference in RMSE appears rather marginal,  
 590 but the associated SAC-SMA parameter values derived with SCE-UA are substantially  
 591 removed from their posterior distribution derived with DREAM. Repeated trials with  
 592 SCE-UA using different values of the algorithmic variables yielded very similar results,  
 593 and did not resolve the problem with premature convergence. We therefore exclude SCE-  
 594 UA from our analysis, and instead consider the SP-UCI method of *Chu et al.* [2010]. This  
 595 new global optimizer was especially designed to overcome some of the main flaws of SCE-  
 596 UA. The results of *Chu et al.* [2010], although limited to a few case studies, demonstrate  
 597 the advantages of SP-UCI over the standard SCE-UA method. *Chu et al.* [2010] kindly  
 598 shared the SP-UCI code with us, and we ran this optimizer sequentially using standard  
 599 settings of the algorithmic variables. Like RWMN and DREAM, each *FE* with SP-UCI  
 600 thus consumes a single *CTU*.

601 Table 4 summarizes the results of the different algorithms used herein. The reported  
 602 statistics denote averages over 25 different trials. We draw two main conclusions from the  
 603 tabulated statistics. In the first place, notice that the three different MCMC methods  
 604 exhibit a very similar effectiveness. The DREAM, DREAM<sub>(zs)</sub> and MT-DREAM<sub>(zs)</sub>  
 605 algorithms consistently converge to the approximate same invariant distribution, with

606 a negligibly small variability among the 25 different trials. A second interesting finding,  
 607 perhaps rather unexpected, is that SP-UCI performs rather poorly with RMSE values that  
 608 are substantially larger than those reported for the different MCMC algorithms. About  
 609 95% of the SP-UCI trials converge prematurely and get stuck in a local basin of attraction  
 610 en route to the global minimum (maximum likelihood) of the posterior distribution. This  
 611 reinforces the severity of the SAC-SMA model calibration problem, and questions the  
 612 ability of SP-UCI to deal with complex and multi-modal search spaces.

613 A third, and final observation is that MT-DREAM<sub>(ZS)</sub> is most efficient in generating  
 614 posterior samples. The  $CTU_{\hat{R}}$ ,  $BURN_{RL}$ , and  $NT_{RL}$  convergence diagnostics demon-  
 615 strate that MT-DREAM<sub>(ZS)</sub> requires the least amount of computational time to explore  
 616 the posterior distribution of the SAC-SMA model parameters. According to the  $CTU_{\hat{R}}$   
 617 criterion, MT-DREAM<sub>(ZS)</sub> is about twice as efficient as the most efficient DREAM<sub>(ZS)</sub>  
 618 parameterization ( $N = 3$ ), and considerably more efficient than the original DREAM al-  
 619 gorithm. This again highlights the advantages of multiple-try sampling. The acceptance  
 620 rate of MT-DREAM<sub>(ZS)</sub> of about 30% is about 3 - 4 times higher than the other adaptive  
 621 Metropolis samplers. This speeds up the efficiency of posterior exploration, and enables  
 622 MT-DREAM<sub>(ZS)</sub> to cope with difficult response surfaces.

623 Although the different diagnostic metrics convey useful information about the perfor-  
 624 mance of the different algorithms, it remains useful to study the convergence behavior of  
 625 the different methods in more detail. Please consider Figure 4 which depicts the evolu-  
 626 tion of the sampled values of the upper zone tension water maximum storage (UZTWM)  
 627 parameter [see, e.g., *Thiemann et al.*, 2001, for details] using the (4a) DREAM, (4b)  
 628 DREAM<sub>(ZS)</sub>, (4c) MT-DREAM<sub>(ZS)</sub>, and (4d) SP-UCI optimization algorithms. Indeed,

the results presented in the various panels are consistent with our previous conclusions.

The MT-DREAM<sub>(ZS)</sub> code converges most rapidly and requires about 2,000 *CTU* to explore the posterior target distribution. This is significantly more efficient than the other two MCMC codes which require about 7,000 (DREAM<sub>(ZS)</sub>) and 50,000 (DREAM) *CTU* to converge to a limiting distribution. Finally, SP-UCI converges rather quickly, but only about 5% of the trials finds the appropriate values of UZTWM.

### 3.4. Groundwater model calibration at the Nevada test site: 241-parameters

The final case study presented herein revisits the work of *Keating et al.* [2010], and involves calibration and predictive uncertainty analysis of a highly parameterized and strongly nonlinear groundwater model. This requires specification of  $d = 241$  different parameters, and constitutes a very difficult optimization problem. A detailed description of the model, site, and data can be found in *Keating et al.* [2010], and so will not be repeated herein. We merely provide a brief synopsis.

Between 1951 and 1992, a total of 659 underground nuclear tests were conducted in Yucca Flat, Nevada Test Site, USA. These explosions have likely enhanced the flux of water into the lower aquifer. A complex three-dimensional hydrogeological flow and transport model was setup to analyze whether the underground tests potentially increased the flux of radionuclides to the groundwater. Yet, this detailed model was shown to be too CPU-intensive to benefit from state-of-the-art optimization and uncertainty analysis methods. A fast-running surrogate model was therefore developed that mimics the key characteristics of the full process model, but is computationally way more efficient.

The design of this surrogate model was based on the simple assumption that for a specific head increase at the location of an explosion, the immediate head perturbation

imposed at a nearby well will be determined by only the distance of the well from the test (i), the time elapsed since the test (ii), and the properties of the rock at the point of measurement (iii). This simplified model contains 241 different parameters, of which 221 are nuclear testing-effect parameters, 10 are rock permeabilities and another 10 are associated with groundwater recharge. As calibration data set we used 361 different head measurements collected at 60 different wells and spanning the time period from 1958 to 2005. In keeping with the previous study, a weighted sum of square residuals (WSSR) was used to measure the distance between the measured head data,  $\hat{y}_i$  and respective predictions,  $y_i(\mathbf{x}, \boldsymbol{\phi})$  of the model:

$$\text{WSSR}(\mathbf{x}|\hat{\mathbf{y}}, \boldsymbol{\phi}) = \sum_{i=1}^{N_m} [w_i (y_i(\mathbf{x}, \boldsymbol{\phi}) - \hat{y}_i)]^2 \quad (12)$$

649 where  $w_i$  are weighting factors [see *Keating et al.*, 2010, for assignment of weights].

650 Previous results presented in *Keating et al.* [2010] have shown that DREAM required  
 651 about 50 million surrogate model evaluations to converge to a limiting parameter distri-  
 652 bution with WSSR values ranging between 370 - 430. This distribution, however cannot  
 653 be considered the actual posterior distribution. In the same paper, it was shown that  
 654 Null Space Monte Carlo [NSMC, *Tonkin and Doherty*, 2009] converged to a very simi-  
 655 lar distribution of parameter values, but found one realization with a significantly better  
 656 WSSR value of about 186 after 402,586 surrogate model evaluations. Finding this solution  
 657 was by no means easy. It required a joint use of manual and computer-based parame-  
 658 ter estimation methods, including the covariance matrix adaptation evolutionary strategy  
 659 [*CMAES, Hansen et al.*, 2003], truncated singular value decomposition [*SVD, Tonkin and*  
 660 *Doherty*, 2005; *Marquardt*, 1963], and automatic user intervention [*AUI, Doherty*, 2009].

<sup>661</sup> We now evaluate the performance of SP-UCI, DREAM<sub>(ZS)</sub>, and MT-DREAM<sub>(ZS)</sub>, and  
<sup>662</sup> test whether they are able to locate the minimum WSSR value of about 186.

Each algorithm was run with a flat prior distribution. The SP-UCI algorithm works directly with the WSSR objective function, but DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub> require specification of a probability density function. We rewrite Eq. (12) in probability space and derive the following log-likelihood function:

$$\pi(\mathbf{x}|\hat{\mathbf{y}}, \boldsymbol{\phi}) \propto -\frac{1}{2} \ln [(2\pi)^d |\boldsymbol{\Sigma}|] - \frac{1}{2} \sum_{i=1}^{N_m} [w_i (y_i(\mathbf{x}, \boldsymbol{\phi}) - \hat{y}_i)]^2 \quad (13)$$

<sup>663</sup> where  $\boldsymbol{\Sigma}$  is a diagonal matrix with nonnegative diagonal elements taken as  $(w_i)^{-2}$ . Note  
<sup>664</sup> that Eq. (13) just provides a different scaling of the WSSR objective function, and thus  
<sup>665</sup> by no means affects the properties of the response surface and location of the posterior  
<sup>666</sup> parameter distribution.

<sup>667</sup> A maximum total of 1,750,000 *CTU* was allowed for the different calibration and pos-  
<sup>668</sup> terior sampling methods using standard settings of the algorithmic variables.

<sup>669</sup> Figure 5 presents the evolution of the sampled WSSR values for (a) SP-UCI (b)  
<sup>670</sup> DREAM<sub>(ZS)</sub> and (c) MT-DREAM<sub>(ZS)</sub>. Each color represents the path of a different trial  
<sup>671</sup> (SP-UCI) or Markov chain (DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>). The triangle, square, and  
<sup>672</sup> cross symbols at the right hand side of the bottom panel report the minimum WSSR values  
<sup>673</sup> found with SP-UCI, DREAM, and PEST, respectively. The advantages of multiple-try  
<sup>674</sup> sampling are immediately obvious. MT-DREAM<sub>(ZS)</sub> has converged to a limiting distribu-  
<sup>675</sup> tion after about 180,000 *CTU* with WSSR values around 200. The minimum WSSR value  
<sup>676</sup> of 193 found with MT-DREAM<sub>(ZS)</sub> is very close to the value of 186 found by PEST. This  
<sup>677</sup> is a remarkable result, and inspires confidence in the ability of MT-DREAM<sub>(ZS)</sub> to solve  
<sup>678</sup> highly parameterized inversion problems. The performance of SP-UCI is rather poor.

679 The code has stagnated after about 50,000 *CTU* to WSSR values of about 910, and is  
 680 unable to escape from this local basin of attraction. Population degeneration has caused  
 681 SP-UCI to prematurely converge. An attempt to introduce particle diversity after about  
 682 150,000 *CTU* (scatter), remains unsuccessful even after 400,000 *CTU* (not shown). Fi-  
 683 nally, DREAM<sub>(ZS)</sub> rapidly converges to WSSR values between 400 - 500, and cannot seem  
 684 to get out of this subspace of solutions. About 1.5 million of additional *CTU* (not shown)  
 685 are necessary for one chain to find WSSR values of around 200, but official convergence  
 686 to the posterior target requires many additional *CTU* (not shown). This finding holds for  
 687 DREAM<sub>(ZS)</sub> parameterized with both  $N = 3$  and  $N = 15$  parallel chains.

688 The results in Fig. 5 not only indicate superior performance of MT-DREAM<sub>(ZS)</sub> but  
 689 also highlight the advantages of stochastic search, and sampling from the past. Whereas  
 690 the population diversity of SP-UCI quickly deteriorates, the Metropolis selection rule of  
 691 MT-DREAM<sub>(ZS)</sub> naturally maintains variability as it is especially designed to converge  
 692 to a distribution of parameter values, rather than a single solution. The various chains  
 693 simulated with MT-DREAM<sub>(ZS)</sub> therefore refuse to settle on a single point, and continue to  
 694 explore a range of WSSR values. This ability to maintain adequate variability is a necessity  
 695 to be able to solve complicated search and optimization problems. The variability sampled  
 696 with DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub> is further enhanced by generating jumps from an  
 697 archive of past solutions.

698 This is further illustrated in Figure 6, which for MT-DREAM<sub>(ZS)</sub> presents the evolution  
 699 of the  $\hat{R}_j$  statistic of *Gelman and Rubin* [1992] for each individual parameter,  $j = 1, \dots, d$ .  
 700 As stated earlier, a value of  $\hat{R}_j < 1.2$  is typically used to declare convergence to a limit-  
 701 ing distribution. Whereas the sampled WSSR values stabilize after about 180,000 *CTU*,

<sup>702</sup> MT-DREAM<sub>(zs)</sub> has not formally converged until approximately 1 million *CTU*. It sim-  
<sup>703</sup> ply takes such a large number of additional samples to adequately explore the  $d = 241$   
<sup>704</sup> dimensional posterior distribution.

<sup>705</sup> The new results with MT-DREAM<sub>(zs)</sub> likely alter some of the previous conclusions of  
<sup>706</sup> *Keating et al.* [2010] who compared formal Bayesian inference with DREAM [*Vrugt et*  
<sup>707</sup> *al.*, 2008] against the more efficient NSMC method of *Tonkin and Doherty* [2009] that  
<sup>708</sup> does not maintain detailed balance and therefore has no theoretical Bayesian foundation.  
<sup>709</sup> Our previous results demonstrated that DREAM and NSCM converged to very similar  
<sup>710</sup> marginal (posterior) parameter distributions. Yet, our results presented herein demon-  
<sup>711</sup> strate substantial differences between the marginal parameter distributions inferred with  
<sup>712</sup> MT-DREAM<sub>(zs)</sub> and their respective pdfs previously derived with DREAM. This is graph-  
<sup>713</sup> ically illustrated in Figure 7 which compares the histogram of the marginal distribution  
<sup>714</sup> of the parameter  $\alpha$  [see *Keating et al.*, 2010] sampled with DREAM (gray bins) and MT-  
<sup>715</sup> DREAM<sub>(zs)</sub> (white bins). Indeed, the marginal distributions deviate considerably from  
<sup>716</sup> each other. This new insight brings into question our previous conclusions, and warrants  
<sup>717</sup> additional analysis.

<sup>718</sup> Altogether, the results in this paper demonstrate some desirable advantages of multi-  
<sup>719</sup> try sampling for estimating complex, multi-modal, and high-dimensional posterior distri-  
<sup>720</sup> butions. The MT-DREAM<sub>(zs)</sub> algorithm not only converges most rapidly from all the  
<sup>721</sup> different MCMC samplers, but also provides samples that most closely represent the un-  
<sup>722</sup> derlying target distribution. This superior performance might be somewhat surprising,  
<sup>723</sup> considering that almost 90% of the samples (with  $k = 5$ ) generated with MT-DREAM<sub>(zs)</sub>  
<sup>724</sup> are thrown away. This appears rather inefficient, yet, the summary statistics presented in

725 the various tables illustrate that multi-try sampling has some desirable advantages. Seem-  
 726 ingly, 1,000 high quality and diverse samples created with MT-DREAM<sub>(ZS)</sub> contain more  
 727 information about the posterior distribution than 10,000 lower quality and less diverse  
 728 samples generated with DREAM<sub>(ZS)</sub> and consisting of many different duplicates.

#### 4. Conclusions

729 Spatially distributed hydrologic models are increasingly being used to study the trans-  
 730 port of water through catchments. These models contain a large number of parameters  
 731 whose values cannot be measured directly in the field but can only meaningfully be ob-  
 732 tained by calibration against a historical record of input - output data. The usefulness  
 733 and applicability of such models thus essentially relies on the availability of powerful  
 734 calibration methods that can efficiently summarize parameter and predictive uncertainty.

735 In this paper, we have presented a novel MCMC algorithm entitled MT-DREAM<sub>(ZS)</sub>  
 736 that is especially designed to solve high-dimensional model calibration problems and sum-  
 737 marize posterior parameter distributions. This method combines the strengths of DREAM  
 738 [Vrugt *et al.*, 2009], sampling from past states [ter Braak and Vrugt, 2008; Vrugt *et al.*,  
 739 2011], snooker updating [ter Braak and Vrugt, 2008], and multiple-try sampling [MTM Liu  
 740 *et al.*, 2000] to evolve an initial population of points to the posterior target distribution.  
 741 MT-DREAM<sub>(ZS)</sub> is especially designed to be implemented on a distributed computing  
 742 cluster. Detailed balance and ergodicity of the algorithm have been studied and ensured,  
 743 which inspires confidence in the ability of MT-DREAM<sub>(ZS)</sub> to generate an example sam-  
 744 ple of the posterior distribution. Four different case studies with different peculiarities  
 745 involving local optima, multi-modality, and high parameter dimensionality, have shown  
 746 that MT-DREAM<sub>(ZS)</sub> is generally superior to existing optimization and search approaches.

747 There are various ways in which the efficiency of MT-DREAM<sub>(ZS)</sub> can be further im-  
 748 proved, particularly for high-dimensional problems. For instance, one could use the ideas  
 749 of *Tonkin and Doherty* [2005] and significantly reduce the dimensionality of the search  
 750 space by reparameterizing the original inverse problem using super and base parameters  
 751 derived from principal component analysis. This could significantly speed up the efficiency  
 752 of posterior sampling, but our initial attempts to date have not been particularly success-  
 753 ful. Not only because each chain works in a different subspace (which makes it difficult  
 754 to define the proposal distribution!), but also because this approach violates detailed bal-  
 755 ance. More generally, coupling MT-DREAM<sub>(ZS)</sub> with dimensionality reduction techniques  
 756 [e.g., *Marzouk and Najm*, 2009] seems attractive to accelerate convergence to a limiting  
 757 distribution. Another idea is to alternate the mix of parallel direction and snooker up-  
 758 dates with a Langevin step [e.g., *Roberts and Rosenthal*, 1998] to sample proposal points  
 759 preferentially in the direction of higher posterior density. This requires explicit knowledge  
 760 of the gradient of  $\pi(\cdot)$ , which can be approximated from the archive of past states. We  
 761 leave these ideas for future work.

762 The source code of MT-DREAM<sub>(ZS)</sub> is written in MATLAB (sequential version) and  
 763 OCTAVE (parallel version) and can be obtained from the second author upon request  
 764 ([jasper@uci.edu](mailto:jasper@uci.edu)).

765 **Acknowledgments.** The authors are thankful to two anonymous referee for their use-  
 766 ful comments and suggestions. We acknowledge Elizabeth Keating for providing the sur-  
 767rogate groundwater model and associated calibration data. We also would like to thank  
 768 Wei Chu for kindly sharing the SP-UCI source code. Computer support, provided by the

<sup>769</sup> SARA center for parallel computing at the University of Amsterdam, The Netherlands,  
<sup>770</sup> is highly appreciated.

## References

- <sup>771</sup> Burnash, R. J. C. (1995), The NWS river forecast systemcatchment modeling, in *Com-*  
<sup>772</sup> *puter Models of Watershed Hydrology*, edited by V. P. Singh, pp. 311–366, Water Re-  
<sup>773</sup> sources Publications, Littleton, CO.
- <sup>774</sup> Chib, S., F. Nardari, and N. Shepard (2002), Markov chain Monte Carlo methods for  
<sup>775</sup> stochastic volatility models, *J. Econometrics*, *108*, 281–316.
- <sup>776</sup> Chu, W., X. Gao, and S. Sorooshian (2010), Improving the shuffled complex evolution  
<sup>777</sup> scheme for optimization of complex nonlinear hydrological systems: Application to the  
<sup>778</sup> calibration of the Sacramento soil moisture accounting model, *Water Resour. Res.*, *46*,  
<sup>779</sup> W09530, doi:10.1029/2010WR009224.
- <sup>780</sup> Dekker, S. C, J.A. Vrugt, R. J. Elkington, (2010), Significant variation in vegetation  
<sup>781</sup> characteristics and dynamics from ecohydrological optimality of net carbon profit. *Eco-*  
<sup>782</sup> *hydrology*, doi:10.1002/eco.177.
- <sup>783</sup> Doherty, J. (2009), PEST: Model independent parameter estimation, Wa-  
<sup>784</sup> termark Numer. Comput., Corinda, Queensland, Australia. (Available at  
<sup>785</sup> <http://www.pesthomepage.org>).
- <sup>786</sup> Duan, Q., V. K. Gupta, and S. Sorooshian (1992), Effective and efficient global optimiza-  
<sup>787</sup> tion for conceptual rainfall-runoff models, *Water. Resour. Res.*, *28*, 1015–1031.
- <sup>788</sup> El Adlouni, S., Favre, A. -C., and Bobée, B. (2006), Comparison of methodologies to  
<sup>789</sup> assess the convergence of Markov chain Monte Marlo methods, *Comput. Stat. Data*  
<sup>790</sup> *An.*, *50*(10), 2685–2701.

- 791 Frenkel, D. (2004), Speed-up of Monte Carlo simulations by sampling of rejected states,  
 792 *Proc. Natl. Acad. Sci. U. S. A.*, 101(51), 457–472, doi:10.703/pnas.0407950101.
- 793 Gelman, A. G., and D. B. Rubin (1992), Inference from iterative simulation using multiple  
 794 sequences, *Stat. Sci.*, 7, 457–472.
- 795 Geweke, J. (1992), Evaluating the accuracy of sampling-based approaches to calculating  
 796 posterior moments, in *Bayesian Statistics 4*, edited by J. M. Bernardo et al., pp. 169–194,  
 797 Clarendon, Oxford, U.K.
- 798 Gilks, W. R., G. O. Roberts, and E. I. George (1994), Adaptive direction sampling,  
 799 *Statistician*, 43, 179–189.
- 800 Gilks, W. R., and G. O. Roberts (1996), Strategies for improving MCMC, in *Markov Chain  
 Monte Carlo in Practice*, edited by W. R. Gilks, S. Richardson, and D. J. Spiegelhalter,  
 801 pp. 89–114, Chapman & Hall, London, U.K.
- 802 Haario, H., E. Saksman, and J. Tamminen (1999), Adaptive proposal distribution for  
 803 random walk Metropolis algorithm, *Comp. Stat.*, 14(3), 375–395.
- 804 Haario, H., E. Saksman, and J. Tamminen (2001), An adaptive Metropolis algorithm,  
 805 *Bernoulli*, 7, 223–242.
- 806 Haario, H., E. Saksman, and J. Tamminen (2005), Componentwise adaptation for high  
 807 dimensional MCMC. *Stat. Comput.*, 20, 265–274.
- 808 Haario H., M. Laine, A. Mira, and E. Saksman (2006), DRAM: Efficient adaptive MCMC.  
 809 *Stat. Comput.*, 16, 339–354.
- 810 Hansen, N., S. D. Muller, and P. Koumoutsakos (2003), Reducing the time complexity  
 811 of the derandomized evolution strategy with covariance matrix adaptation (CMAES),  
 812 *Evol. Comput.*, 11(1), 1–18, doi:10.1162/106365603321828970.

- 814 Hastings, H. (1970), Monte Carlo sampling methods using Markov chains and their ap-  
 815 plications, *Biometrika*, 57, 97–109.
- 816 Keating, E. H., J. Doherty, J. A. Vrugt, and Q. Kang (2010), Optimization and un-  
 817 certainty assessment of strongly nonlinear groundwater models with high parameter  
 818 dimensionality, *Water. Resour. Res.*, 46, W10517, doi:10.1029/2009WR008584.
- 819 Kuczera, G., D. Kavetski, B. Renard, and M. Thyre (2010), A limited memory accelera-  
 820 tion strategy for MCMC sampling in hierarchical Bayesian calibration of hydrological  
 821 models, *Water. Resour. Res.*, 46, W07602, doi:10.1029/2009WR008985.
- 822 Laloy, E., Fasbender D., and C.L. Bielders (2010a), Parameter optimization and uncer-  
 823 tainty analysis for plot-scale continuous modeling of runoff using a formal Bayesian  
 824 approach, *J. Hydrol.*, 380(1-2), 82-93, doi:10.1016/j.jhydrol.2009.10.025.
- 825 Laloy, E., Weynants M., C.L. Bielders, M. Vanclooster, and M. Javaux (2010b), How  
 826 efficient are one-dimensional models to reproduce the hydrodynamic behavior of struc-  
 827 tured soils subjected to multi-step outflow experiments?, *J. Hydrol.*, 393(1-2), 37-52,  
 828 doi:10.1016/j.jhydrol.2010.02.017.
- 829 Liu, J. S., F. Liang, and W. H. Wong (2000), The multiple-try method and lo-  
 830 cal optimization in metropolis sampling, *J. Am. Stat. Assoc.*, 95(449), 121–134,  
 831 doi:10.2307/2669532.
- 832 Marquardt, D. W. (1963) An algorithm for least-squares estimation of nonlinear param-  
 833 eters, *J. Soc. Indust. Appl. Math.*, 11, 431–441.
- 834 Marzouk, Y. M, and H. N Najm (2009), Dimensionality reduction and polynomial chaos  
 835 acceleration of Bayesian inference in inverse problems, *J. Comput. Phys.*, 118.6, 862–  
 836 1902, doi:10.1016/j.jcp.2008.11.024.

- 837 Metropolis, N., A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller (1953),  
 838     Equation of state calculations by fast computing machines, *J. Chem. Phys.*, **21**, 1087–  
 839     1092.
- 840 Murray, I. (2007), Advances in Markov chain Monte Carlo methods, *Ph.D. Thesis*, Uni-  
 841     versity of London.
- 842 Radu, V. C., J. Rosenthal, and C. Yang (2009), Learn from the thy neighbor: parallel-  
 843     chain and regional adaptive MCMC, *J. Am. Stat. Assoc.*, **104**(488), 1454–1466.
- 844 Raftery, A.E., and S.M. Lewis (1992), One long run with diagnostics: Implementation  
 845     strategies for Markov chain Monte Carlo. *Stat. Sci.*, **7**, 493–497.
- 846 Roberts, G. O., and W. R. Gilks (1994), Convergence of adaptive direction sampling, *J.*  
 847     *Multivariate Anal.*, **49**, 287–298.
- 848 Roberts, G. O., A. Gelman, and W. R. Gilks (1997), Weak convergence and optimal  
 849     scaling of random walk Metropolis algorithms. *Ann. Appl. Prob.*, **7**, 110–120.
- 850 Roberts, G. O., and J.S. Rosenthal (1998), Optimal scaling of discrete approximations to  
 851     Langevin diffusions, *J. R. Statist. Soc. B*, **60**, 255–268.
- 852 Roberts, G. O., and J. S. Rosenthal (2007), Coupling and ergodicity of adaptive Markov  
 853     chain Monte Carlo algorithms, *J. Appl. Probab.*, **44**, 458–475.
- 854 Scharnagl, B., J.A. Vrugt, H. Vereecken, and M. Herbst (2010), Information content of  
 855     incubation experiments for inverse modeling of carbon pools in the Rothamsted model:  
 856     a Bayesian approach, *Biogeosciences*, **7**, 763–776.
- 857 Schoups, G., and J. A. Vrugt (2010), A Formal Likelihood Function for Parameter and  
 858     Predictive Inference of Hydrologic Models with Correlated, Heteroscedastic and Non-  
 859     Gaussian Errors, *Water. Resour. Res.*, **46**, W10531, doi:10.1029/2009WR008933.

- 860 ter Braak, C. J. F. (2006), A Markov Chain Monte Carlo version of the genetic algo-  
 rithm differential evolution: Easy Bayesian computing for real parameter space, *Stat.*  
*Comput.*, 16(3), 239249, doi:10.1007/s11222-006-8769-1.
- 863 ter Braak, C. J. F., and J. A. Vrugt (2008), Differential evolution Markov chain with  
 864 snooker updater and fewer chains, *Stat. Comput.*, 18(4), 435446, doi:10.1007/s11222-  
 865 008-9104-9.
- 866 Thiemann, M., M. Trossset, H. Gupta, and S. Sorooshian (2001), Bayesian recursive  
 867 parameter estimation for hydrologic models, *Water. Resour. Res.*, 37(10), 25212535,  
 868 doi:10.1029/2000WR900405.
- 869 Tonkin, M., and J. Doherty (2005), A hybrid regularized inversion methodology  
 870 for highly parameterized environmental models, *Water. Resour. Res.*, 41, W10412,  
 871 doi:10.1029/2005WR003995.
- 872 Tonkin, M., and J. Doherty (2009), Calibration constrained Monte Carloanalysis of highly  
 873 parameterized models using subspace techniques, *Water. Resour. Res.*, 45, W00B10,  
 874 doi:10.1029/2007WR006678.
- 875 Vrugt, J. A., H. V. Gupta, W. Bouten, and S. Sorooshian (2003), A Shuffled  
 876 Complex Evolution Metropolis algorithm for optimization and uncertainty assess-  
 877 ment of hydrologic model parameters, *Water. Resour. Res.*, 39(8), art. No. 1201,  
 878 doi:10.1029/2002WR001642.
- 879 Vrugt, J. A., H. V. Gupta, S. C. Dekker, S. Sorooshian, T. Wagener, and W. Bouten  
 880 (2006), Application of stochastic parameter optimization to the Sacramento soil mois-  
 881 ture accounting model, *J. Hydrol.*, 325(14), 288–307.

- 882 Vrugt, J. A., C. J. F. ter Braak, M. P. Clark, J. M. Hyman, and B. A. Robinson  
883 (2008), Treatment of input uncertainty in hydrologic modeling: doing hydrology back-  
884 ward with Markov chain Monte Carlo simulation, *Water. Resour. Res.*, 44, W00B09,  
885 doi:10.1029/2007WR006720.
- 886 Vrugt, J. A., C. J. F. ter Braak, C. G. H. Diks, D. Higdon, B. A. Robinson, and J. M.  
887 Hyman (2009), Accelerating Markov chain Monte Carlo simulation by differential evo-  
888 lution with self-adaptive randomized subspace sampling, *Int. J. Nonlinear Sci. Numer.*  
889 *Simul.*, 10(3), 273–290.
- 890 Vrugt, J. A., E. Laloy, C. J. F. ter Braak, and J. M. Hyman (2011), Posterior exploration  
891 using differential evolution adaptive Metropolis with sampling from past states, *to be*  
892 *submitted to SIAM*.
- 893 Yapo, P. O., H. V. Gupta, and S. Sorooshian (1996), Calibration of conceptual rainfall-  
894 runoff models: sensitivity to calibration data. *J. Hydrol.*, 181, 23–48.

**Table 1.** Performance of MT-DREAM<sub>(ZS)</sub> against RWMN, DREAM, and DREAM<sub>(ZS)</sub> for the 200-dimensional Gaussian distribution with correlated dimensions, using a maximum total of 1,000,000 (RWMN, DREAM) and 400,000 (DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>) Computational Time Units (*CTU*). *N* is the number of chains, *D* measures closeness to the true posterior target,  $CTU_{\hat{R}}$ ,  $BURN_{RL}$  and  $NT_{RL}$  are the *Gelman and Rubin* [1992] and the two *Raftery and Lewis* [1992] convergence criteria expressed in computational time, respectively, *IF* is the inefficiency factor and *AR* is the average acceptance rate. *D* and *IF* are computed using the last 250,000 draws generated by each method within the allowed computational time. Reported values represent averages over 25 independent runs.

	<i>N</i>	<i>D</i>	$CTU_{\hat{R}}$ $\times 10^4$ <i>CTU</i>	$BURN_{RL}$ $\times 10^3$ <i>CTU</i>	$NT_{RL}$ $\times 10^6$ <i>CTU</i>	<i>IF</i> [-]	<i>AR</i> %
RWMN	1	0.524	<i>N/A</i> *	30.29	10.541	807.5	24.3
DREAM	200	0.086	<i>N/C</i> **	0.43	0.681	6.3	17.0
DREAM <sub>(ZS)</sub>	3	0.062	3.975	0.21	0.447	121.0	16.8
DREAM <sub>(ZS)</sub>	15	0.062	1.602	0.03	0.059	114.7	17.3
MT-DREAM <sub>(ZS)</sub>	3	0.038	2.959	0.35	0.239	53.1	45.2

\*Not applicable.

\*\*None of the 25 runs have converged within the allowed  $10^6$  *CTU*.

**Table 2.** Performance of MT-DREAM<sub>(ZS)</sub> against RWMN, DREAM, and DREAM<sub>(ZS)</sub> for the 25-dimensional tri-modal distribution and a maximum total of 2,000,000 (RWMN, DREAM) and 400,000 (DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>) Computational Time Units (*CTU*). *N* is the number of chains, *D* measures closeness to the true posterior target,  $CTU_{\hat{R}}$ ,  $BURN_{RL}$  and  $NT_{RL}$  are the *Gelman and Rubin* [1992] and the two *Raftery and Lewis* [1992] convergence criteria expressed in computational time, respectively, *IF* is the inefficiency factor and *AR* is the average acceptance rate. *D* and *IF* are computed using the last 500,000 draws generated by each method within the allowed computational time. Reported values represent averages over 25 independent runs.

	<i>N</i>	<i>D</i>	$CTU_{\hat{R}}$ $\times 10^4$ <i>CTU</i>	$BURN_{RL}$ $\times 10^3$ <i>CTU</i>	$NT_{RL}$ $\times 10^6$ <i>CTU</i>	<i>IF</i> [-]	<i>AR</i> %
RWMN	1	0.830	<i>N/A</i> *	0.24	0.503	8.1	24.1
DREAM	25	0.087	113.041	16.24	32.652	73.4	5.9
DREAM <sub>(ZS)</sub>	5	0.191	3.640	7.25	12.841	120.9	9.8
DREAM <sub>(ZS)</sub>	25	0.085	7.836	0.56	0.894	50.1	9.8
MT-DREAM <sub>(ZS)</sub>	5	0.052	3.360	0.85	1.768	196.1	28.6

\*Not applicable.

**Table 3.** Description of the SAC-SMA model parameters, including their prior and 95% posterior uncertainty intervals derived with MT-DREAM<sub>(ZS)</sub>.

Parameter	Units	Prior	Posterior
UZTWM	mm	1.0 - 150.0	12.31 - 29.34
UZFWM	mm	1.0 - 150.0	10.14 - 37.30
LZTWM	mm	1.0 - 500.0	211.76 - 255.17
LZFPM	mm	1.0 - 1000.0	99.38 - 146.29
LZFSM	mm	1.0 - 1000.0	37.30 - 84.48
ADIMP	[-]	0.0 - 0.40	0.29 - 0.38
UZK	day <sup>-1</sup>	0.1 - 0.5	0.15 - 0.49
LZPK	day <sup>-1</sup>	0.0001 0.025	0.0101 - 0.0194
LZSK	day <sup>-1</sup>	0.01 0.25	0.23 - 0.25
ZPERC	[-]	1.0 250.0	155.45 - 249.32
REXP	[-]	1.0 5.0	2.29 - 4.42
PCTIM	[-]	0.0 0.1	0.0001 - 0.0063
PFREE	[-]	0.0 0.6	0.0002 - 0.1477

**Table 4.** Comparison of MT-DREAM<sub>(zs)</sub> against, DREAM, DREAM<sub>(zs)</sub>, and SP-UCI for the SAC-SMA model calibration and a maximum total of 100,000 Computational Time Units ( $CTU$ ). Listed statistics represent averages from 25 different trials. The variables  $N$ ,  $CTU_{\hat{R}}$ ,  $BURN_{RL}$ ,  $NT_{RL}$ ,  $IF$  and  $AR$  have been defined in the main text and previous table captions.  $IF$  is computed using the last 25,000 draws generated by each method within the allowed computational time.  $E_{MEAN}$ ,  $E_{MIN}$ , and  $E_{MAX}$  denote the average, minimum, and maximum best RMSE of the 25 trials. The MCMC standard error being in the range  $0.04 - 0.07 \text{ m}^3 \text{ s}^{-1}$  for the Metropolis samplers, RMSE values are rounded to the first decimal. The variable  $N_{LOC}$  reports the number of runs being stuck into a local minima ( $RMSE > 13.3 \text{ m}^3 \text{ s}^{-1}$ ).

	$N$	$E_{MEAN}$	$E_{MIN}$	$E_{MAX}$	$N_{LOC}$	$CTU_{\hat{R}}$	$BURN_{RL}$	$NT_{RL}$	$IF$	$AR$
		$\text{m}^3 \text{ s}^{-1}$	$\text{m}^3 \text{ s}^{-1}$	$\text{m}^3 \text{ s}^{-1}$		$\times 10^4 \text{ CTU}$	$\times 10^3 \text{ CTU}$	$\times 10^6 \text{ CTU}$	[ ]	%
DREAM	13	13.2	13.1	13.2	0	8.586*	99.18	90.253	48.8	6.1
DREAM <sub>(zs)</sub>	3	13.2	13.1	13.2	0	0.712	2.70	4.710	172.4	7.3
DREAM <sub>(zs)</sub>	15	13.1	13.1	13.2	0	1.553	3.74	5.927	62.5	5.8
MT-DREAM <sub>(zs)</sub>	3	13.1	13.1	13.2	0	0.427	2.22	3.521	52.9	24.9
SP-UCI		<i>N/A</i> **	13.6	13.1	14.0	24			<i>N/A</i> **	

\*12 from the 25 runs did not appropriately converge within the allowed  $10^5 \text{ CTU}$ .

\*\*Not applicable.

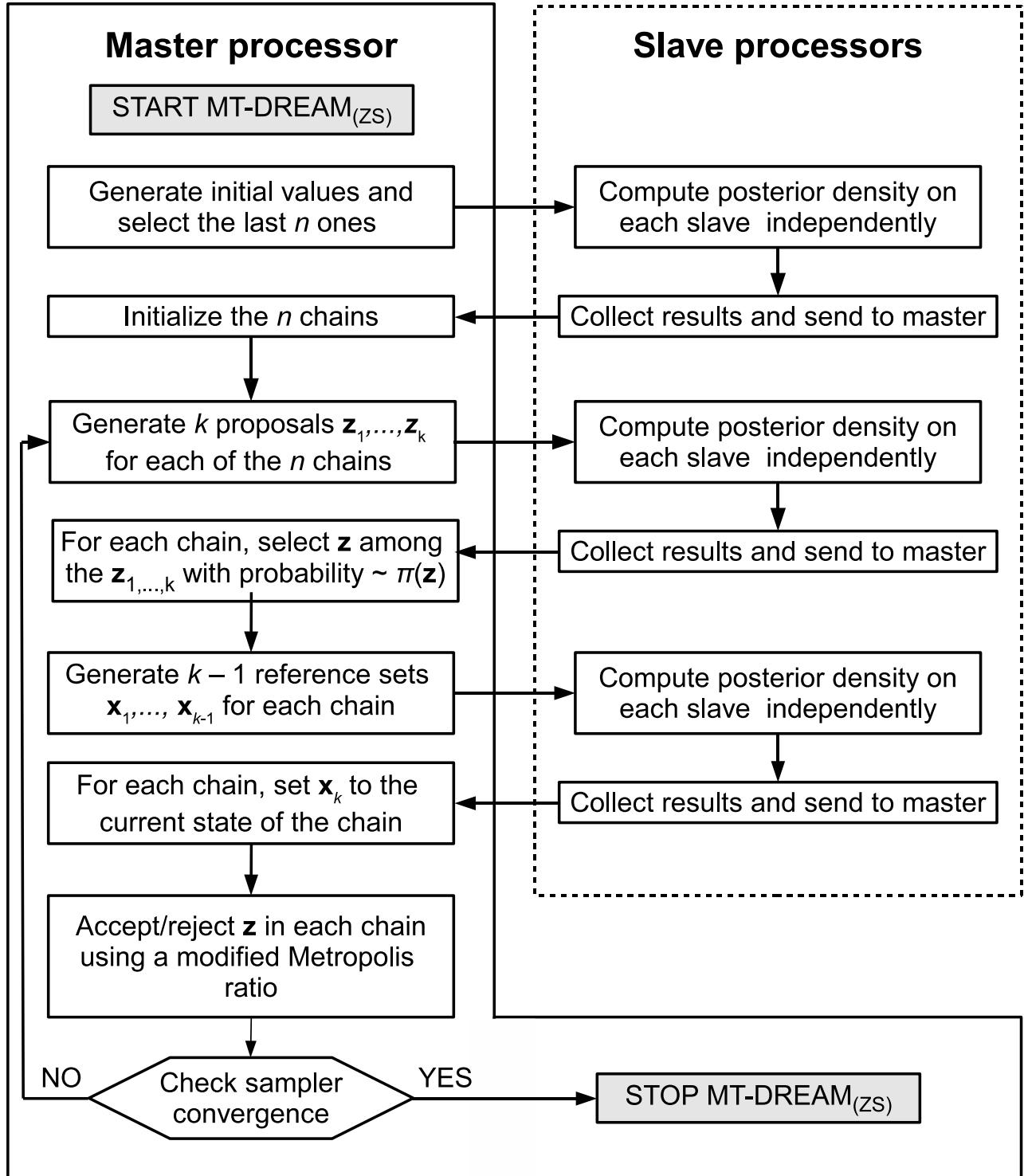
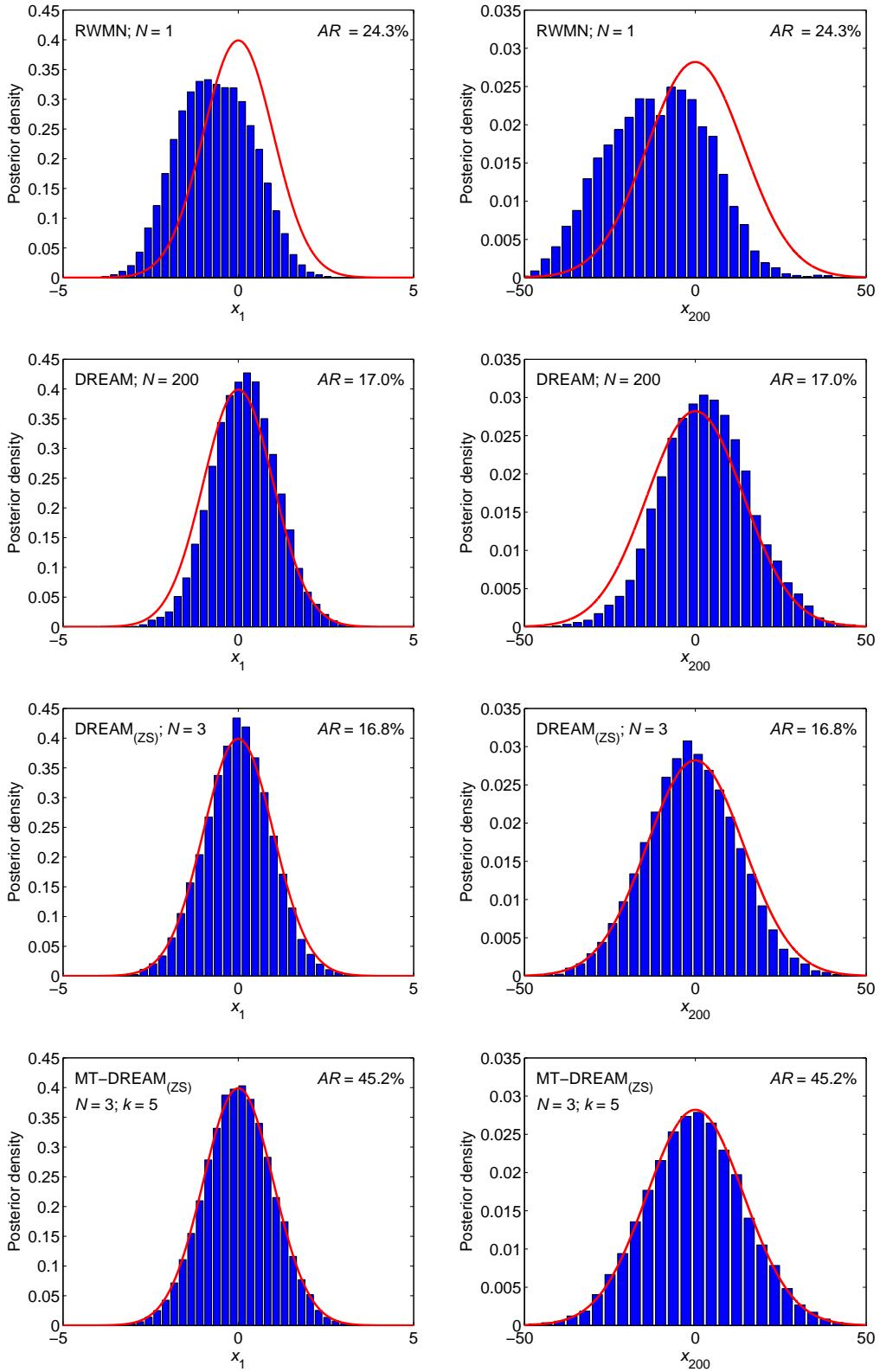
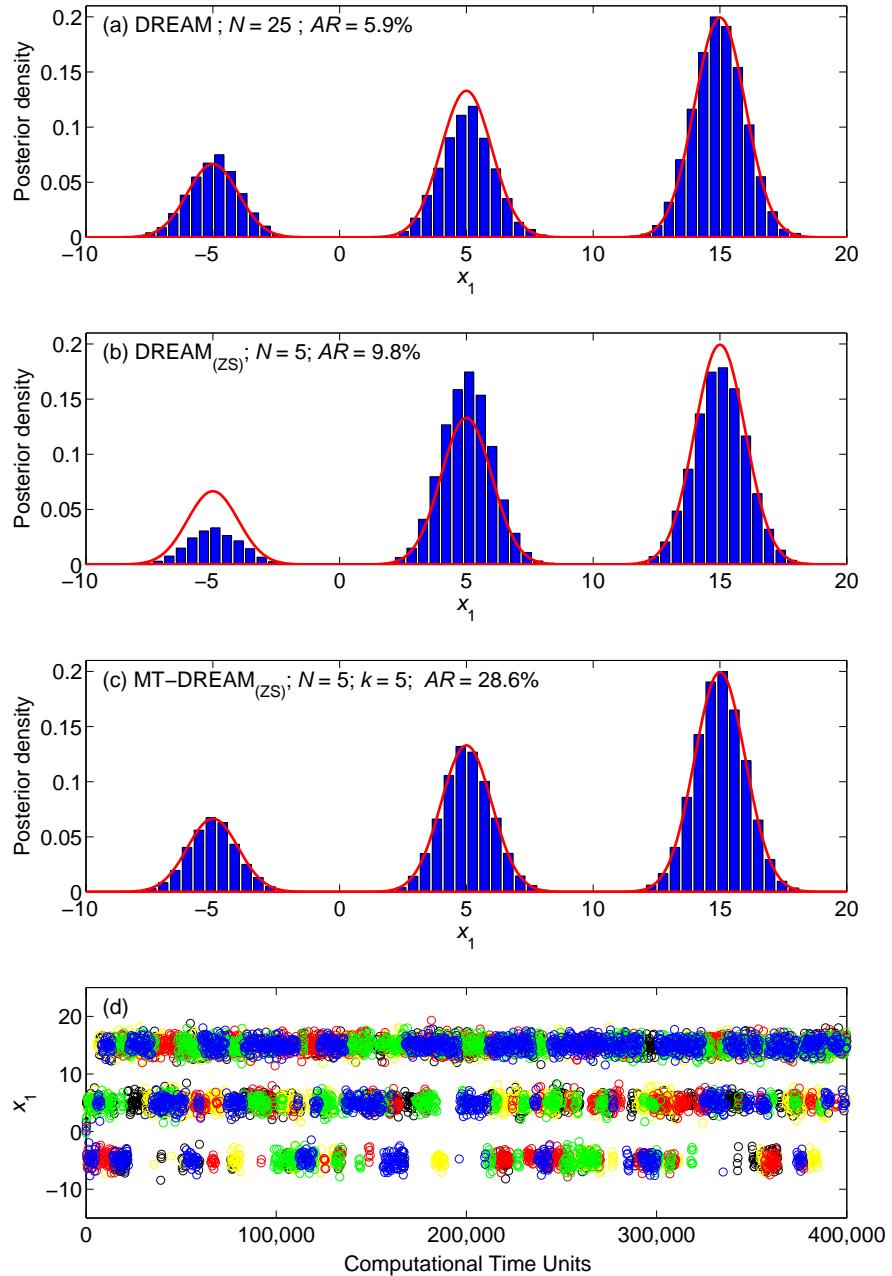


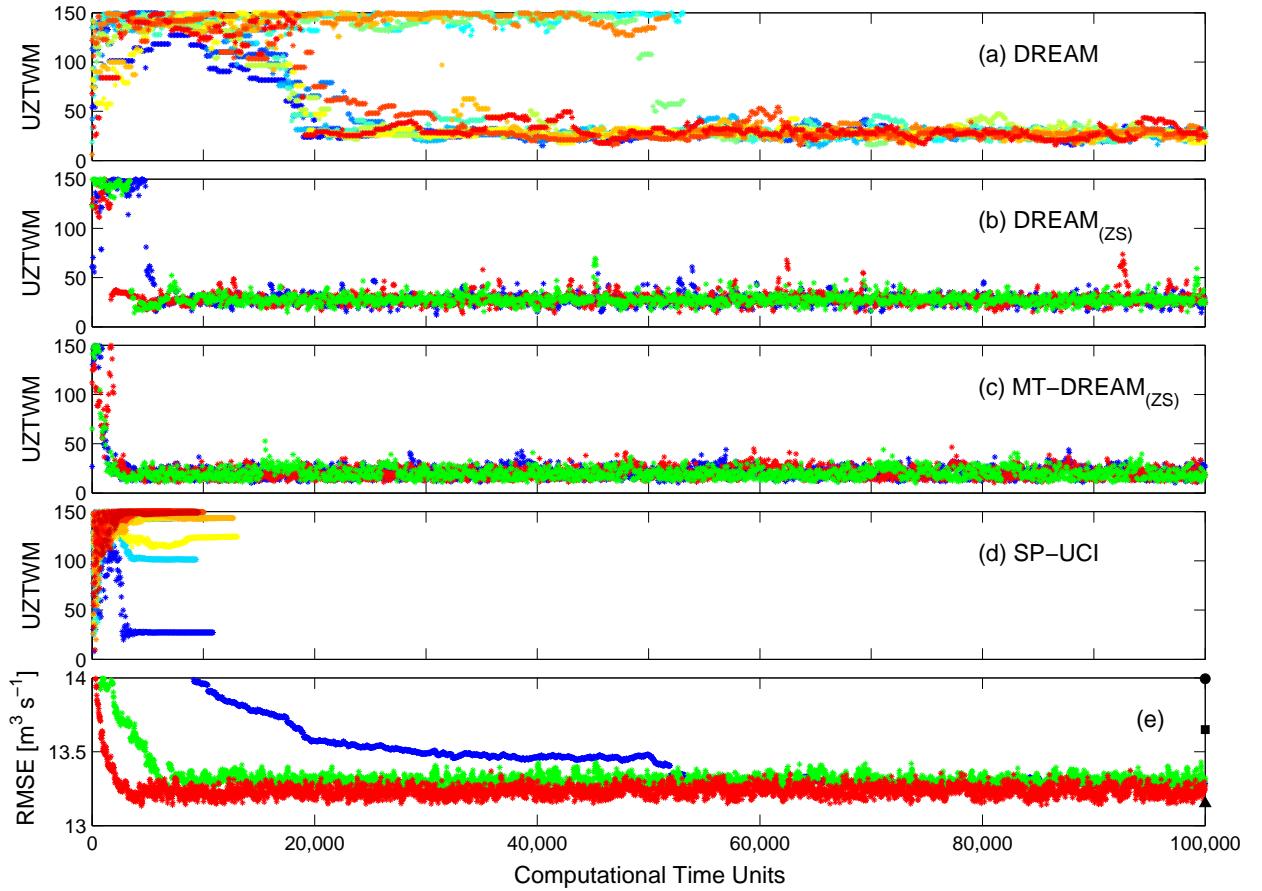
Figure 1. Sketch of the distributed MT-DREAM<sub>(ZS)</sub> algorithm.



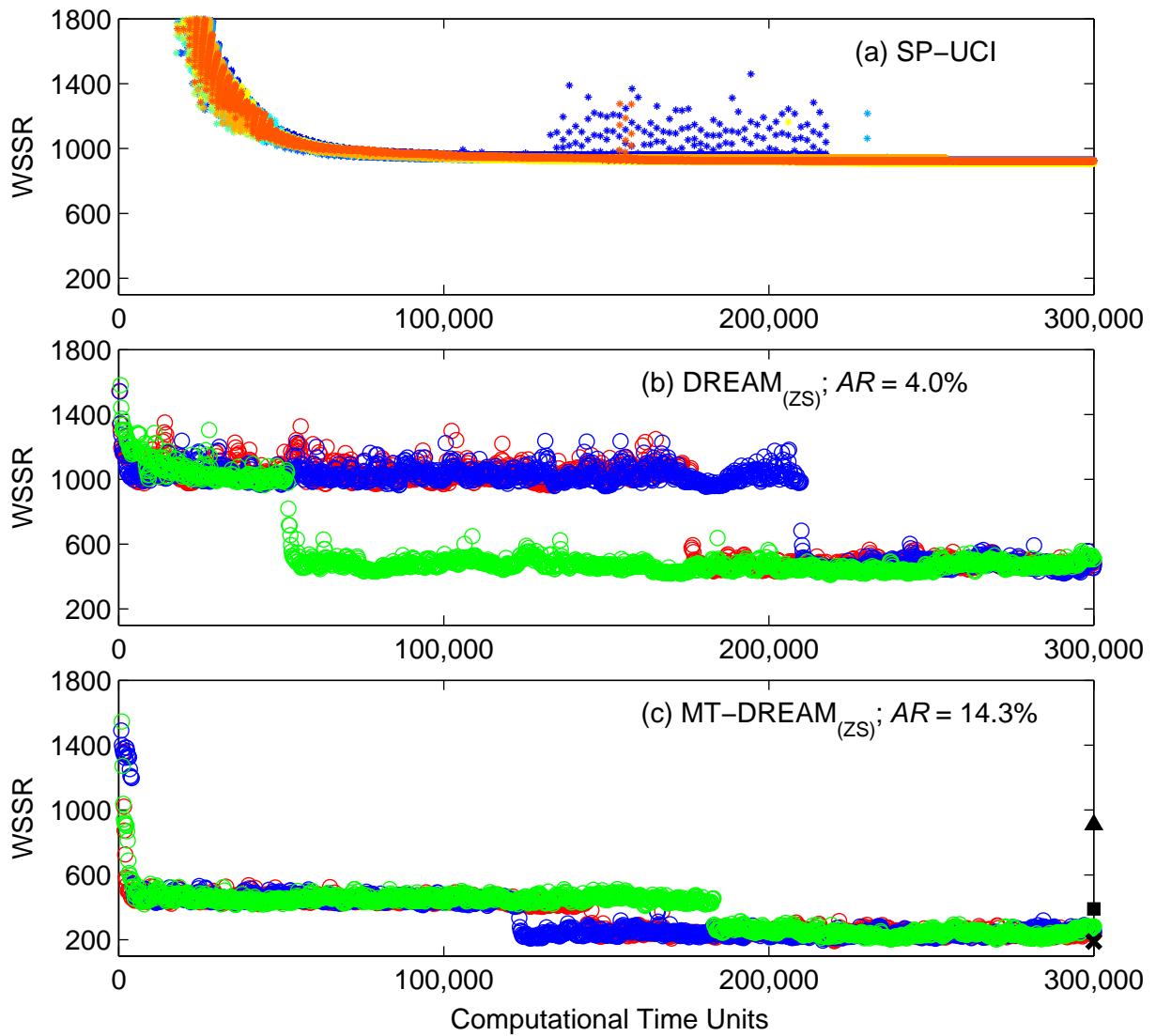
**Figure 2.** Marginal posterior pdfs of dimensions 1 ( $x_1$ ; left column) and 200 ( $x_{200}$ ; right column) for the 200-dimensional multivariate normal distribution with correlated dimensions using RMWN (top panel), DREAM (upper middle panel), DREAM<sub>(ZS)</sub> (lower middle panel), and MT-DREAM<sub>(ZS)</sub> (bottom panel). Each individual plot reports  $N$ , the number of chains, and  $AR$ , the average acceptance rate. For MT-DREAM<sub>(ZS)</sub>, we also list the value of  $k$ , the number of parallel proposal points. The true marginal distributions are indicated with a solid red line.



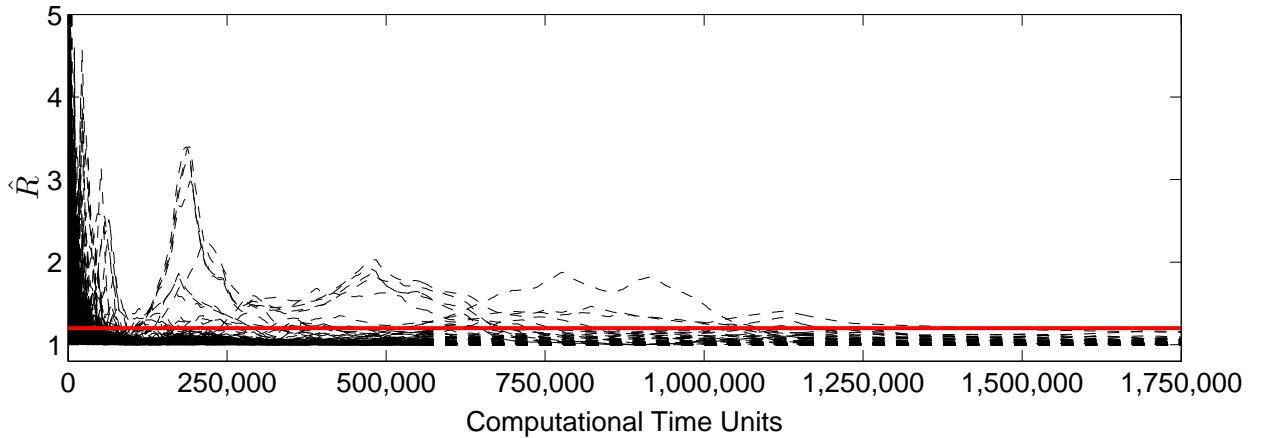
**Figure 3.** Marginal posterior pdfs of dimension 1 ( $x_1$ ) obtained with (a) DREAM , (b) DREAM<sub>(ZS)</sub> and (c) MT-DREAM<sub>(ZS)</sub>. The variable  $N$  is the number of chains, and  $AR$  denotes the average acceptance rate. For MT-DREAM<sub>(ZS)</sub>, the value of  $k$ , the number of parallel proposal points, is also listed. The histograms are derived from the last 50% of the samples in the joint chains. The true posterior pdf is indicated with the solid red line. The bottom graph depicts the evolution of the  $N$  pathways sampled with MT-DREAM<sub>(ZS)</sub>. These results demonstrate a superior performance of MT-DREAM<sub>(ZS)</sub>.



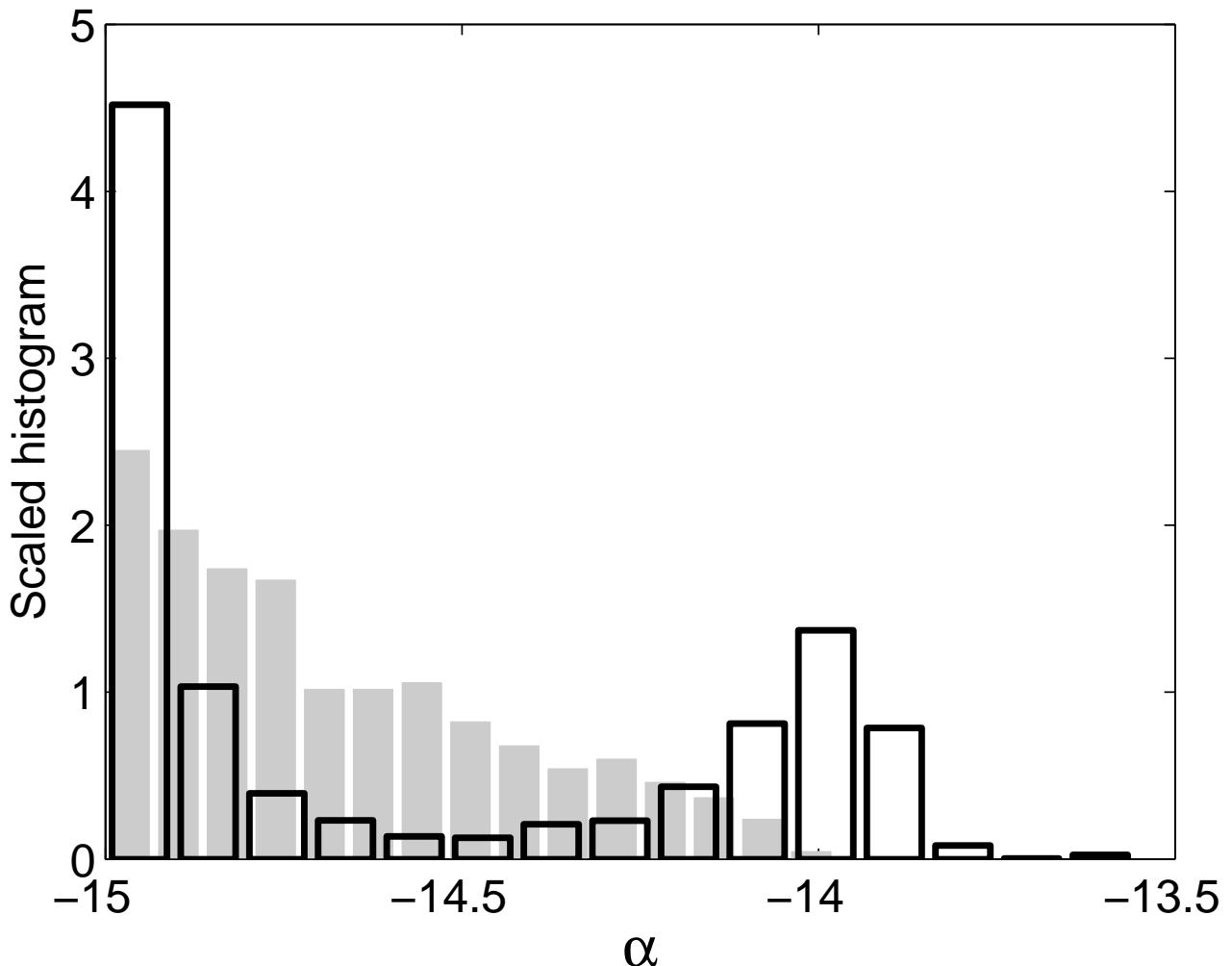
**Figure 4.** Evolution of sampled values of the upper zone tension water maximum storage (UZTWM) parameter (in mm) with the (a) DREAM, (b) DREAM<sub>(ZS)</sub> with  $N = 3$ , (c) MT-DREAM<sub>(ZS)</sub> MCMC sampling schemes, and (d) SP-UCI global optimization algorithm. Each chain in panels (a)-(c) is coded with a different color. For SP-UCI, we use color coding to illustrate each of the 25 different runs. The bottom panel (e) depicts the evolution of the mean RMSE value of the  $N$  different chains derived with DREAM (blue), DREAM<sub>(ZS)</sub> with  $N = 3$  (green), and MT-DREAM<sub>(ZS)</sub> (red). The black square, black dot and black triangle in the bottom panel (e) represent the mean, maximum and minimum RMSE value of the 25 different SP-UCI runs, respectively.



**Figure 5.** Trace plots of the sampled WSSR values using (a) SP-UCI, (b) DREAM<sub>(ZS)</sub> and (c) MT-DREAM<sub>(ZS)</sub> for the 241-dimensional groundwater model calibration problem. The variable  $AR$  has been defined in the main text. Colors are used to indicate different trials (SP-UCI) or Markov chains (DREAM<sub>(ZS)</sub> and MT-DREAM<sub>(ZS)</sub>). The different symbols at the right hand side of the bottom panel summarize the minimum WSSR values found with SP-UCI (triangle), DREAM (square), and PEST (cross), respectively.



**Figure 6.** MT-DREAM<sub>(ZS)</sub>: Evolution of the  $\hat{R}$ -statistic of *Gelman and Rubin* [1992] during the MT-DREAM<sub>(ZS)</sub> run for the 241-dimensional groundwater model calibration problem of [*Keating et al.*, 2010]. The different parameters  $\hat{R}_j, j = 1, \dots, d$  are coded with dashed lines. If the dashed-lines jointly fall below the solid red line, convergence to a limiting distribution can be officially declared.



**Figure 7.** Histograms of the marginal posterior pdf of the  $\alpha$  parameter in the groundwater model of Keating *et al.* [2010] derived with DREAM (gray bins) and MT-DREAM<sub>(ZS)</sub> (white bins). The marginal posterior distributions differ quite substantially, questioning some of the main conclusions of our previous work.