

¹ **Approximate Bayesian Computation Using Markov**
² **Chain Monte Carlo Simulation: DREAM_(ABC)**

Mojtaba Sadegh,¹ and Jasper A. Vrugt^{1,2,3}

¹Department of Civil and Environmental
Engineering, University of California, Irvine,
CA, USA.

²Department of Earth System Science,
University of California, Irvine, CA, USA.

³IBG-3 Agrosphere, Forschungszentrum
Julich, Julich, Germany.

3 **Abstract.** The quest for a more powerful method for model evaluation
4 has inspired *Vrugt and Sadegh* [2013] to introduce "likelihood-free" inference
5 as vehicle for diagnostic model evaluation. This class of methods is also re-
6 ferred to as Approximate Bayesian Computation (ABC) and relaxes the need
7 for a residual-based likelihood function in favor of one or multiple different
8 summary statistics that exhibit superior diagnostic power. Here, we propose
9 several methodological improvements over commonly used ABC sampling
10 methods to permit inference of complex system models. Our methodology
11 entitled, DREAM_(ABC) uses the DiffeRential Evolution Adaptive Metropo-
12 lis algorithm [*Vrugt et al.*, 2008, 2009] as its main building block and takes
13 advantage of a continuous fitness function to efficiently explore the behav-
14 ioral model space. Three case studies demonstrate that DREAM_(ABC) is at
15 least 3 - 1,000 times more efficient than commonly used ABC sampling meth-
16 ods.

1. Introduction and Scope

Bayesian methods have become increasingly popular for fitting hydrologic models to data (e.g. streamflow, water chemistry, groundwater table depth, soil moisture, pressure head, snow water equivalent). Bayes' rule updates the prior probability of a certain hypothesis when new data, $\tilde{\mathbf{Y}} = \{\tilde{y}_1, \dots, \tilde{y}_n\}$ (also referred to as evidence) become available. The hypothesis typically constitutes the parameter values, $\boldsymbol{\theta}$, of a model, F , which simulates the observed data using

$$\tilde{\mathbf{Y}} \leftarrow F(\boldsymbol{\theta}, \tilde{\mathbf{u}}, \tilde{\mathbf{x}}_0) + \mathbf{e}, \quad (1)$$

where $\tilde{\mathbf{u}} = \{\tilde{u}_1, \dots, \tilde{u}_n\}$ denotes the observed forcing data, $\tilde{\mathbf{x}}_0$ signifies the initial state of the system, and $\mathbf{e} = \{e_1, \dots, e_n\}$ includes observation error, as well as error due to the fact that the simulator, F , may be systematically different from the real system of interest, $\mathfrak{S}(\boldsymbol{\theta})$, for the parameters $\boldsymbol{\theta}$. If our main interest is in the parameters of the model, Bayes law is given by

$$p(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) = \frac{p(\boldsymbol{\theta})p(\tilde{\mathbf{Y}}|\boldsymbol{\theta})}{p(\tilde{\mathbf{Y}})}, \quad (2)$$

where $p(\boldsymbol{\theta})$ denotes the prior parameter distribution, $p(\tilde{\mathbf{Y}}|\boldsymbol{\theta}) \equiv L(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$ is the likelihood function, and $p(\tilde{\mathbf{Y}})$ represents the evidence. As all statistical inferences of the parameters can be made from the unnormalized density, we conveniently remove $p(\tilde{\mathbf{Y}})$ from the denominator and write $p(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$.

The likelihood function, $L(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$, summarizes, in probabilistic sense, the overall distance between the model simulation and corresponding observations. The mathematical definition of this function has been subject to considerable debate in the hydrologic and

35 statistical literature [e.g. *Schoups and Vrugt* [2010]; *Smith et al.* [2010]; *Evin et al.* [2013]].
36 Simple likelihood functions that assume Gaussian error residuals are statistically conve-
37 nient, but this assumption is often not borne out of the probabilistic properties of the
38 error residuals that show significant variations in bias, variance, and autocorrelation at
39 different parts of the simulated watershed response. Such non-traditional residual distri-
40 butions are often caused by forcing data and model structural errors, whose probabilistic
41 properties are very difficult, if not impossible, to adequately characterize. This makes it
42 rather difficult, if not impossible, to isolate and detect epistemic errors (model structural
43 deficiencies), a prerequisite to improving our understanding and theory of water flow and
44 storage in watersheds.

45 The inability of classical likelihood-based fitting methods to detect model malfunction-
46 ing is evident if we critically assess the progress that has been made in modeling of the
47 rainfall-runoff transformation. For instance, consider the Sacramento soil moisture ac-
48 counting (SAC-SMA) model introduced by *Burnash et al.* [1973] in the early 1970s and
49 used by the US National Weather Service for flash-flood forecasting throughout the United
50 States. In about four decades of fitting the SAC-SMA model to (spatially distributed)
51 streamflow data, we have not been able to make any noticeable improvements to the
52 underlying equations of the model. This is even more disturbing given the relative low
53 order complexity of the SAC-SMA model. If for such relatively simple (lumped) hydro-
54 logic models our fitting methods are unable to illuminate to what degree a representation
55 of the real world has been adequately achieved and how the model should be improved,
56 the prospects of learning and scientific discovery for the emerging generation of very high
57 order system models are rather poor, because more complex process representations lead

58 (unavoidably) to greater interaction among model components, and perpetually larger
59 volumes of field and remote sensing data need to be utilized for system characterization
60 and evaluation.

61 The limitations of classical residual-based fitting methods has stimulated *Gupta et al.*
62 [2008] (amongst others) to propose a signature-based approach to model evaluation. By
63 choosing the signatures so that they each measure different but relevant parts of system
64 behavior, diagnostic evaluation proceeds with analysis of the behavioural similarities (and
65 differences) of the observed data and corresponding model simulations. Ideally, these
66 differences are then related to individual process descriptions, and correction takes place
67 by refining/improving these respective components of the model. What is left is the
68 numerical implementation of diagnostic model evaluation.

69 In a previous paper, *Vrugt and Sadegh* [2013] advocated the use of "likelihood-free"
70 inference for diagnostic model evaluation. This approach, introduced in the statistical
71 literature about three decades ago [*Diggle and Gratton*, 1984], is especially useful for
72 cases where the likelihood is intractable, too expensive to be evaluated, or impossible
73 to be formulated explicitly. This class of methods is also referred to as Approximate
74 Bayesian Computation (ABC), a term coined by *Beaumont et al.* [2002], and widens the
75 realm of models for which statistical inference can be considered [*Marjoram et al.*, 2003;
76 *Sisson et al.*, 2007; *Del Moral et al.*, 2008; *Joyce and Marjoram*, 2008; *Grelaud et et al.*,
77 2009; *Ratmann et al.*, 2009]. ABC has rapidly gained popularity in the past few years,
78 in particular for the analysis of complex problems arising in population genetics, ecology,
79 epidemiology, and systems biology. The first application of ABC in hydrology can be
80 found in *Nott et al.* [2012] and establishes a theoretical connection between ABC and

81 GLUE-based approaches. Other work on this topic can be found in the recent publication
 82 by *Sadegh and Vrugt* [2013].

83 The premise behind ABC is that $\boldsymbol{\theta}^*$ should be a sample from the posterior distribution if
 84 the distance between the observed and simulated data, $\rho(\tilde{\mathbf{Y}}, \mathbf{Y}(\boldsymbol{\theta}^*))$, is smaller than some
 85 small positive value, ϵ [*Marjoram et al.*, 2003; *Sisson et al.*, 2007]. Figure 1 provides a
 86 conceptual overview of the ABC methodology. All ABC based methods approximate the
 87 likelihood function by simulations, the outcomes of which are compared with the observed
 88 data [*Beaumont*, 2010; *Bertorelle et al.*, 2010; *Csilléry et al.*, 2010]. In so doing, ABC
 89 algorithms attempt to approximate the posterior distribution by sampling from

$$p(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto \int_{\mathcal{Y}} p(\boldsymbol{\theta}) \text{Model}(\mathbf{y}|\boldsymbol{\theta}) \mathbf{I}(\rho(\tilde{\mathbf{Y}}, \mathbf{Y}(\boldsymbol{\theta})) \leq \epsilon) d\mathbf{y}, \quad (3)$$

90 where \mathcal{Y} denotes the support of the simulated data, $\mathbf{Y} \sim \text{Model}(\mathbf{y}|\boldsymbol{\theta})$, and $\mathbf{I}(a)$ is an
 91 indicator function that returns one if the condition a is satisfied and zero otherwise. The
 92 accuracy of the estimated posterior distribution, $p(\boldsymbol{\theta}|\rho(\tilde{\mathbf{Y}}, \mathbf{Y}(\boldsymbol{\theta})) \leq \epsilon)$ depends on the value
 93 of ϵ . In the limit of $\epsilon \rightarrow 0$ the sampled distribution will converge to the true posterior,
 94 $p(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$ [*Pritchard et al.*, 1999; *Beaumont et al.*, 2002; *Ratmann et al.*, 2009; *Turner and*
 95 *van Zandt*, 2012]. Yet, this requires the underlying model operator to be stochastic, and
 96 hence $\text{Model}(\cdot)$ in Equation (3) is equivalent to the output of the deterministic model F
 97 in Equation (1) plus a random error with probabilistic properties equal to those of \mathbf{e} .

98 For sufficiently complex system models and/or large data sets, it will be difficult, if not
 99 impossible, to find a model simulation that always fits the data within ϵ . It has therefore
 100 become common practice in ABC to use one or more summary statistics of the data rather
 101 than the data itself. Ideally, these chosen summary statistics, $S(\cdot)$ are sufficient and thus

102 provide as much information for the model parameters as the original data set itself. In
103 practice, however, the use of summary statistics usually entails a loss of information and
104 hence results in an approximate likelihood, especially for complex models. Partial least
105 squares [Wegmann *et al.*, 2009] and information-theory [Barnes *et al.*, 2011] can help
106 to determine (approximately) a set of nearly-sufficient marginal statistics. Nonetheless,
107 complex models admitting sufficient statistics are practical exceptions.

108 The most common ABC algorithm implements simple rejection sampling which relies
109 on satisfying the condition $\rho(S((\tilde{\mathbf{Y}})), S(\mathbf{Y}(\boldsymbol{\theta}^*))) \leq \epsilon$. This method has the practical ad-
110 vantage of being relatively easy to implement and use, but its efficiency depends critically
111 on the choice of the prior sampling distribution. If this prior distribution is a poor approx-
112 imation of the actual posterior distribution, then many of the proposed samples will be
113 rejected. This leads to dramatically low acceptance rates, and thus excessive CPU-times.
114 Indeed, *Vrugt and Sadegh* [2013] and *Sadegh and Vrugt* [2013] report acceptance rates of
115 less than 0.1% for hydrologic models with just a handful of parameters. One remedy to
116 this problem is to increase the value of ϵ , but this leads to an inaccurate approximation
117 of the posterior distribution.

118 A number of methodological advances have been proposed to enhance the sampling
119 efficiency of ABC algorithms. One common approach is to use a set of monotonically
120 decreasing ϵ values. This allows the algorithm to sequentially adapt the prior distribution
121 and converge to a computationally feasible final value of ϵ . Nonetheless, these algorithms
122 still rely on a boxcar kernel (step function) to evaluate the fitness of each sample, and
123 are not particularly efficient in high dimensional search spaces. In this paper we intro-
124 duce a Markov Chain Monte Carlo (MCMC) simulation method that enhances, some-

125 times dramatically, the ABC sampling efficiency. This general-purpose method entitled,
126 DREAM_(ABC) uses the DiffeRential Evolution Adaptive Metropolis algorithm [*Vrugt et*
127 *al.*, 2008, 2009] as its main building block, and replaces the indicator function in Equation
128 (3) with a continuous kernel to decide whether to accept candidate points or not. The
129 proposed methodology is benchmarked using synthetic and real-world simulation experi-
130 ments.

131 The remainder of this paper is organized as follows. In section 2 we summarize the re-
132 sults of commonly used ABC sampling methods by application to a synthetic benchmark
133 study. Section 3 introduces the main elements of the DREAM_(ABC) algorithm and dis-
134 cusses several of its advantages. This is followed in Section 4 with two synthetic and one
135 real-world simulation experiment. In this section we are especially concerned with sam-
136 pling efficiency and robustness. Finally, section 5 concludes this paper with a discussion
137 and summary of our main findings.

2. Approximate Bayesian Computation

138 The ABC method provides an excellent vehicle for diagnostic model evaluation [*Vrugt*
139 *and Sadegh*, 2013] by using one or multiple different summary statistics that, when rooted
140 in the relevant environmental theory, should have a much stronger and compelling diag-
141 nostic power than some residual-based likelihood function. Challenges lie in the proper
142 selection of summary metrics that adequately extract all the available information from
143 the calibration data, how to deal with input data uncertainty, how to detect epistemic
144 errors (lack of knowledge), how to determine an appropriate (small) value for ϵ , and how
145 to efficiently sample complex multi-dimensional spaces involving many tens to hundreds of
146 parameters. This paper is focused on the last topic, and proposes several methodological

147 developments to overcome the shortcomings of standard ABC sampling methods. The
 148 other topics will be investigated in subsequent papers.

149 We first discuss two common ABC sampling methods that have found widespread
 150 application and use within the context of likelihood-free inference. We then introduce
 151 DREAM_(ABC), a Markov chain Monte Carlo (MCMC) implementation of ABC that per-
 152 mits inference of complex system models.

2.1. Rejection Algorithm

153 Once the summary statistic(s) has(have) been defined we are left with finding all those
 154 values of θ^* for which $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\theta^*))) \leq \epsilon$. The most basic algorithm to do so uses
 155 rejection sampling. This algorithm proceeds as follows

Algorithm 1 ABC-Rejection Sampler

```

1: for  $i = 1, \dots, N$  do
2:   while  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\theta^*))) > \epsilon$  do
3:     Sample  $\theta^*$  from the prior,  $\theta^* \sim p(\theta)$ 
4:     Simulate data  $\mathbf{Y}$  using  $\theta^*$ ,  $\mathbf{Y} \sim \text{Model}(\theta^*)$ 
5:     Calculate  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\theta^*)))$ 
6:   end while
7:   Set  $\theta^i \leftarrow \theta^*$ 
8:   Set  $w^i \leftarrow \frac{1}{N}$ 
9: end for

```

156 In words, the ABC rejection (ABC-REJ) algorithm proceeds as follows. First we sample
 157 a candidate point, θ^* , from some prior distribution, $p(\theta)$. We then use this proposal to
 158 simulate the output of the model, $\mathbf{Y} \sim \text{Model}(\theta^*)$. We then compare the simulated data,
 159 \mathbf{Y} , with the observed data, $\tilde{\mathbf{Y}}$, using a distance function, $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\theta^*)))$. If this dis-
 160 tance function is smaller than some small positive tolerance value, ϵ then the simulation
 161 is close enough to the observations that the candidate point, θ^* has some nonzero proba-

162 bility of being in the approximate posterior distribution, $\hat{p}(\boldsymbol{\theta}|\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta})))) \leq \epsilon$. By
 163 repeating this process N times, ABC-REJ provides an estimate of the actual posterior
 164 distribution.

165 Unfortunately, standard rejection sampling method typically requires massive computa-
 166 tional resources to generate a sufficient number of samples from the posterior distribution.
 167 Failure to maintain an adequate sampling density may result in under sampling probable
 168 regions of the parameter space. This inefficiency can provide misleading results, par-
 169 ticularly if $p(\boldsymbol{\theta}|\tilde{\mathbf{Y}})$ is high dimensional and occupies only a small region interior to the
 170 prior distribution. Only if $p(\boldsymbol{\theta})$ is a good representation of the actual posterior parameter
 171 distribution then ABC-REJ can achieve an adequate sampling efficiency.

2.2. Population Monte Carlo Simulation

172 To guarantee convergence to the appropriate limiting distribution, the value of ϵ in
 173 Algorithm 1 (ABC-REJ) needs to be taken very small. Values of $0.01 \leq \epsilon \leq 0.05$ are
 174 often deemed appropriate. Unfortunately, this will produce very low acceptance rates,
 175 particularly if the prior distribution is poorly chosen and extends far beyond the posterior
 176 distribution. To increase sampling efficiency, it would seem logical to stepwise reduce the
 177 value of ϵ and to use the accepted samples to iteratively adapt the prior distribution. This
 178 is the principal idea behind population Monte Carlo (PMC) algorithms. These methods
 179 are used extensively in physics and statistics for many-body problems, lattice spin systems
 180 and Bayesian inference, and also referred to as "quantum Monte Carlo", "transfer-matrix
 181 Monte Carlo", "Monte Carlo filter", "particle filter" and "sequential Monte Carlo". The
 182 PMC sampler of *Beaumont et al.* [2009]; *Turner and van Zandt* [2012] is specifically
 183 designed for ABC-inference and works as follows

Algorithm 2 ABC-Population Monte Carlo sampler

```

1: At iteration  $j = 1$ ,
2: for  $i = 1, \dots, N$  do
3:   while  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*))) > \epsilon_1$  do
4:     Sample  $\boldsymbol{\theta}^*$  from the prior,  $\boldsymbol{\theta}^* \sim p(\boldsymbol{\theta})$ 
5:     Simulate data  $\mathbf{Y}$  using  $\boldsymbol{\theta}^*$ ,  $\mathbf{Y} \sim \text{Model}(\boldsymbol{\theta}^*)$ 
6:     Calculate  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*)))$ 
7:   end while
8:   Set  $\boldsymbol{\Theta}_1^i \leftarrow \boldsymbol{\theta}^*$ 
9:   Set  $\mathbf{w}_1^i \leftarrow \frac{1}{N}$ 
10: end for
11: Set  $\boldsymbol{\Sigma}_1 \leftarrow 2\text{Cov}(\boldsymbol{\Theta}_1)$ ,
12: At iteration  $j > 1$ ,
13: for  $j = 2, \dots, J$  do
14:   for  $i = 1, \dots, N$  do
15:     while  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^{**}))) > \epsilon_j$  do
16:       Sample  $\boldsymbol{\theta}^*$  from the previous population,  $\boldsymbol{\theta}^* \sim \boldsymbol{\Theta}_{j-1}$  with probability  $\mathbf{w}_{j-1}$ 
17:       Perturb  $\boldsymbol{\theta}^*$  by sampling  $\boldsymbol{\theta}^{**} \sim \mathcal{N}_d(\boldsymbol{\theta}^*, \boldsymbol{\Sigma}_{j-1})$ 
18:       Simulate data  $\mathbf{Y}$  using  $\boldsymbol{\theta}^{**}$ ,  $\mathbf{Y} \sim \text{Model}(\boldsymbol{\theta}^{**})$ 
19:       Calculate  $\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^{**})))$ 
20:     end while
21:     Set  $\boldsymbol{\Theta}_j^i \leftarrow \boldsymbol{\theta}^{**}$ 
22:     Set  $w_j^i \leftarrow \frac{p(\boldsymbol{\theta}_j^i)}{\sum_{u=1}^N \mathbf{w}_{j-1}^u q_d(\boldsymbol{\theta}_{j-1}^u | \boldsymbol{\theta}_j^i, \boldsymbol{\Sigma}_{j-1})}$ 
23:   end for
24:   Set  $\boldsymbol{\Sigma}_j \leftarrow 2\text{Cov}(\boldsymbol{\Theta}_j)$ 
25: end for

```

(4)

184 In short, the ABC-PMC sampler starts out as ABC-REJ during the first iteration, $j = 1$,
 185 but using a much larger initial value for ϵ . This will significantly enhance the initial ac-
 186 ceptance rate. During each successive iteration, $j = \{2, \dots, J\}$, the value of ϵ is decreased
 187 and the multi-normal proposal distribution, $q = \mathcal{N}_d(\boldsymbol{\theta}_{j-1}^k, \boldsymbol{\Sigma}_{j-1})$, adapted using $\boldsymbol{\Sigma}_{j-1} =$
 188 $2\text{Cov}(\boldsymbol{\theta}_{j-1}^1, \dots, \boldsymbol{\theta}_{j-1}^N)$ with $\boldsymbol{\theta}_{j-1}^k$ drawn from a multinomial distribution, $\mathfrak{F}(\boldsymbol{\Theta}_{j-1} | \mathbf{w}_{j-1})$,
 189 where $\boldsymbol{\Theta}_{j-1} = \{\boldsymbol{\theta}_{j-1}^1, \dots, \boldsymbol{\theta}_{j-1}^N\}$ is a $N \times d$ matrix and $\mathbf{w}_{j-1} = \{w_{j-1}^1, \dots, w_{j-1}^N\}$ is a N -
 190 vector of normalized weights, $\sum_{i=1}^N w_{j-1}^i = 1$ and $w_{j-1}^i \geq 0$. In summary, a sequence of
 191 (multi)normal proposal distributions is used to iteratively refine the samples and explore
 192 the posterior distribution. This approach, similar in spirit as the adaptive Metropolis
 193 sampler of [Haario *et al.*, 1999, 2001], achieves a much higher sampling efficiency than
 194 ABC-REJ, particularly for cases where the prior distribution, $p(\boldsymbol{\theta})$, is a poor approxima-
 195 tion of the actual target distribution.

196 The PMC sampler of *Turner and van Zandt* [2012] assumes that the sequence of ϵ
 197 values is specified by the user. Practical experience suggests that a poor selection of
 198 $\boldsymbol{\epsilon} = \{\epsilon_1, \dots, \epsilon_J\}$ can lead to very low acceptance rates or even premature convergence.
 199 *Sadegh and Vrugt* [2013] have therefore introduced an alternative variant of ABC-PMC
 200 with adaptive selection of $\epsilon_{j(j>1)}$. This method requires the user to specify only the initial
 201 kernel bandwidth, ϵ_1 , and subsequent values of $\boldsymbol{\epsilon}$ are determined from the $\rho(\cdot)$ values of
 202 the N most recently accepted samples. This approach is not only more practical, but
 203 also enhances convergence speed to the posterior distribution. We therefore prescribe the
 204 sequence of $\boldsymbol{\epsilon}$ values in ABC-PMC using the outcome of several adaptive runs.

205 It is interesting to note that the PMC sampler has elements in common with genetic
 206 algorithms (GA) [Higuchi, 1997] in that a population of individuals is used to search the

parameter space. The main difference between both approaches is that PMC is specifically designed for statistical inference of the marginal and joint parameter distributions, whereas GAs are specialized in optimization. Yet, it is not difficult to modify the PMC sampler so that it converges to a single "best" solution. Nonetheless, one should be particularly careful using common GA operators such as crossover and mutation. Such genetic operators can significantly improve the search capabilities of ABC-PMC in high dimensional search spaces, but can harm convergence properties.

To benchmark the efficiency of ABC-REJ and ABC-PMC, we start by fitting a relatively simple mixture of two Gaussian distributions, which has become a classical problem in the ABC literature [Sisson *et al.*, 2007; Beaumont *et al.*, 2009; Toni *et al.*, 2009; Turner and Sederberg, 2012]

$$p(\theta) = \frac{1}{2}\mathcal{N}\left(0, \frac{1}{100}\right) + \frac{1}{2}\mathcal{N}(0, 1), \quad (5)$$

where $\mathcal{N}(a, b)$ is a normal distribution with mean, a and standard deviation, b . The solid black line in Figure 2 plots the actual target distribution.

We now test the efficiency of ABC-REJ and ABC-PMC using the following distance function,

$$\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*))) = \begin{cases} |\bar{\mathbf{Y}}| & \text{with probability of 50\%} \\ |y_1| & \text{with probability of 50\%} \end{cases}$$

where $\mathbf{Y} = \{y_1, \dots, y_{100}\}$, $y_i \sim \mathcal{N}(\theta, 1)$, and the operator $|\cdot|$ signifies the modulus (absolute value). In keeping with the statistical literature, we assume a uniform prior distribution, $p(\theta) \sim \mathcal{U}[-10, 10]$, with $\epsilon = 0.025$ (ABC-REJ) and $\epsilon = \{1, 0.75, 0.5, 0.25, 0.1, 0.05, 0.025\}$ (ABC-PMC). Fig. 2 plots histograms of the ABC-REJ (A: left plot) and ABC-PMC (B: middle plot) derived posterior distribution of θ using $N = 1,000$ samples.

227 The marginal distribution derived with ABC-REJ and ABC-PMC are in good agreement
228 with the known target distribution (black line). Table 1 summarizes the performance of
229 the ABC-REJ and ABC-PMC sampler. The rejection sampler (ABC-REJ) requires about
230 386,742 function evaluations to find $N = 1,000$ behavioral solutions. This corresponds
231 to an acceptance rate (AR, %) of approximate 0.26%, which can be considered highly
232 inefficient. The ABC-PMC sampler on the other hand requires fewer function evaluations
233 (170,848) to explore the target distribution, with an acceptance rate of about 0.59%. Note
234 however that ABC-PMC underestimates the sampling density of points in the tails of the
235 posterior suggesting that the algorithm has not been able to fully explore $p(\theta)$. The results
236 of DREAM_(ABC) that are listed at the bottom of Table 1 will be discussed in section 4.1
237 of this paper.

3. Markov Chain Monte Carlo Simulation

238 The adaptive capabilities of the ABC-PMC sampler offer significant computational ad-
239 vantages over ABC-REJ. However, further methodological improvements are warranted
240 to enable inference of complex simulation models involving high dimensional parameter
241 spaces. The use of a boxcar fitness kernel (zero probability everywhere except for a small
242 interval where it is a constant) is theoretically convenient, but makes it very difficult for
243 any sampling algorithm to determine the preferred search direction. All rejected simula-
244 tions receive a similar score, irrespective of whether their $\rho(\cdot)$ values are in close proximity
245 of the threshold, ϵ or far removed. This is certainly not desirable and unnecessarily com-
246 plicates posterior exploration. Furthermore, ABC-REJ and ABC-PMC update all entries
247 of the parameter vector simultaneously. This is equivalent to a crossover of 100%, and
248 adequate for low-dimensional problems involving just a handful of parameters, but not

249 necessarily efficient in high dimensional search spaces. For such problems, conditional (or
 250 subspace) sampling has desirable properties and can enhance, sometimes dramatically,
 251 the speed of convergence. The method we propose herein is more commonly known as
 252 Metropolis-within-Gibbs, and samples individual dimensions (or groups of parameters) in
 253 turn.

3.1. Continuous fitness kernel

254 A boxcar kernel has the disadvantage that all samples with $\rho(\cdot)$ value larger than ϵ
 255 are considered equal and discarded. This is the basis of rejection sampling, and can
 256 be very inefficient particularly if the prior sampling distribution is poorly chosen. To
 257 solve this problem, *Turner and Sederberg* [2012] recently introduced an alternative ABC
 258 sampling approach using MCMC simulation with a continuous fitness kernel. This method
 259 is based on the concept of noisy-ABC [*Beaumont et al.*, 2002; *Blum and François*, 2010]
 260 and perturbs the model simulation with a random error, ξ

$$\mathbf{Y} \leftarrow \text{Model}(\boldsymbol{\theta}^*) + \boldsymbol{\xi} \quad (6)$$

261 If we assume that $\boldsymbol{\xi}$ follows a multivariate normal distribution, $\mathcal{N}(\mathbf{0}_n, \alpha)$, then we can
 262 evaluate the probability density, $p(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*)))$, of $\boldsymbol{\theta}^*$ using

$$p(\boldsymbol{\theta}^*|\alpha) = \frac{1}{\alpha\sqrt{2\pi}} \exp\left(-\frac{1}{2}\alpha^{-2}\left(\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*)))\right)^2\right), \quad (7)$$

263 where α is an algorithmic (or free) parameter that defines the width of the kernel. As a
 264 result, the approximation in Equation (3) becomes [*Turner and Sederberg*, 2012]

$$p(\boldsymbol{\theta}|\tilde{\mathbf{Y}}) \propto \int_{\mathbf{y}} p(\boldsymbol{\theta}) \text{Model}(\mathbf{y}|\boldsymbol{\theta}) p\left(\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*)))\right) d\mathbf{y} \quad (8)$$

265 This approach has recently been coined kernel-based ABC (KABC) in the statistical lit-
 266 erature, and opens up an arsenal of advanced Monte Carlo based sampling methods to
 267 explore the posterior distribution. Three preliminary case studies of different complexity
 268 (not shown herein) demonstrate that KABC with DREAM is at least 3 – 1,000 times
 269 more efficient than ABC-PMC. Unfortunately, KABC can run into a fatal problem. The
 270 MCMC algorithm may produce a nice bell shaped posterior distribution but with simu-
 271 lated summary statistics that are far removed from their observed values. For example,
 272 let's assume an extreme case in which the model is unable to fit any of the observed
 273 summary statistics within the required tolerance ϵ . Sampling would still produce a limit-
 274 ing distribution, but with probability densities of Equation (7) that are practically zero.
 275 These results are undesirable, and have nothing to do with the actual MCMC method
 276 used, replacing this with another sampling approach would give similar findings. The cul-
 277 prit is the continuous kernel of Equation (7), which does not a-priori bound the feasible
 278 space of the posterior solutions. Changing Equation (7) to a boxcar function with $p(\cdot) = 1$
 279 in the interval $[-\epsilon, \epsilon]$ around the measured summary statistic(s), and exponential decline
 280 of the density outside this interval runs into the same problems and is thus also futile.

281 We here propose an alternative method for fitness assignment that helps convergence
 282 of a MCMC simulator to the posterior distribution. We define the fitness of $\boldsymbol{\theta}^*$ as follows

$$f(\boldsymbol{\theta}^*, \phi) = \phi - \rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*))), \quad (9)$$

283 where $\phi > 0$ is a coefficient that bounds the fitness values between $(-\infty, \phi]$. The smaller
 284 the distance of the simulated summary metrics to their observed values, the higher the
 285 fitness. By setting $\phi = \epsilon$, then $f(\boldsymbol{\theta}^*) \in [0, \epsilon]$ is a necessary condition for a sample, $\boldsymbol{\theta}^*$, to be

286 called a posterior solution, otherwise the sample is non behavioral and can be discarded.
 287 This condition is easily verified a-posteriori from the sampled fitness values of the Markov
 288 chains.

289 We are now left with a definition of the selection rule to help determine whether to
 290 accept trial moves or not. The original scheme proposed by *Metropolis et al.* [1953] was
 291 constructed using the condition of detailed balance. If $p(u)$ ($p(i)$) denotes the probability
 292 to find the system in state u (i) and $q(u \rightarrow i)$ ($q(i \rightarrow u)$) is the conditional probability to
 293 perform a trial move from u to i (i to u), then the probability $p_{\text{acc}}(u \rightarrow i)$ to accept the
 294 trial move from u to i is related to $p_{\text{acc}}(i \rightarrow u)$ according to:

$$p(u)q(u \rightarrow i)p_{\text{acc}}(u \rightarrow i) = p(i)q(i \rightarrow u)p_{\text{acc}}(i \rightarrow u) \quad (10)$$

295 If we assume a symmetric jumping distribution, that is $q(u \rightarrow i) = q(i \rightarrow u)$, then it
 296 follows that

$$\frac{p_{\text{acc}}(u \rightarrow i)}{p_{\text{acc}}(i \rightarrow u)} = \frac{p(i)}{p(u)} \quad (11)$$

297 This Equation does not yet determine the acceptance probability. *Metropolis et al.* [1953]
 298 made the following choice:

$$p_{\text{acc}}(u \rightarrow i) = \min \left[1, \frac{p(i)}{p(u)} \right], \quad (12)$$

299 to determine whether to accept a trial move or not. This selection rule has become
 300 the basic building block of many existing MCMC algorithms. *Hastings* [1970] extended
 301 Equation (12) to non-symmetrical jumping distributions in which $q(u \rightarrow i) \neq q(i \rightarrow u)$.

302 Unfortunately, Equation (9) is not a proper probability density function, and hence ap-
 303 plication of (12) will lead to a spurious approximation of the posterior distribution. The

304 same problem arises with other fitness functions derived from Equation (7), for exam-
 305 ple, the Nash-Sutcliffe efficiency, pseudolikelihoods and a composite of multiple objective
 306 functions. We therefore calculate the acceptance probability using

$$p_{\text{acc}}(u \rightarrow i) = \max\left(\mathbb{I}(f(i) \geq f(u)), \mathbb{I}(f(i) \geq (\phi - \epsilon))\right) \quad (13)$$

307 where $\mathbb{I}(\cdot)$ is an indicator function, and the operator $\max(\mathbb{I}(a), \mathbb{I}(b))$ returns one if a
 308 and/or b is true, and zero otherwise. Thus, $p_{\text{acc}}(u \rightarrow i) = 1$ (we accept) if the fitness of
 309 the proposal i is higher than or equal to that of the current state of the chain, u . On the
 310 contrary, if the fitness of i is smaller than that of u then $p_{\text{acc}}(u \rightarrow i) = 0$ and the proposal
 311 is rejected, unless $f(i) \geq 0$, then we still accept.

312 The binary acceptance probability of Equation (13) differs fundamentally from a regular
 313 MCMC selection rule, but has important practical advantages that promulgate converge to
 314 the correct limiting distribution. Initially, when the samples are rather inferior, Equation
 315 (13) enforces the MCMC algorithm to act as an optimizer and only accept proposals
 316 with a higher fitness and thus smaller distance to the observed values of the summary
 317 statistics. The transitions of the Markov chain during this time are irreversible with a
 318 zero acceptance probability of the previous state (backward jump). This changes the
 319 moment a candidate point has been sampled whose fitness, $f(\cdot) \geq 0$. From this point
 320 forward, the acceptance probability of Equation (13) leads to a reversible Markov chain
 321 and the successive samples can be used to approximate the posterior target distribution.

322 An important limitation of Equation (13) is that it cannot incorporate non-symmetric
 323 jump distributions, such as the snooker updater used in DREAM_(ZS) and MT-DREAM_(ZS)

324 [*Laloy and Vrugt, 2012*]. This would require a Hastings-type correction, but cannot be
 325 readily incorporated within the current framework.

3.2. Pseudo-code of DREAM_(ABC)

326 We can solve for the posterior distribution of Equation (9) using MCMC simulation
 327 with DREAM [*Vrugt et al., 2008, 2009*]. This method uses subspace sampling to increase
 328 search efficiency and overcome some of the main limitations of ABC-REJ and ABC-PMC.
 329 In DREAM, K ($K > 2$) different Markov chains are run simultaneously in parallel, and
 330 multivariate proposals are generated on the fly from the collection of chains, Θ^{t-1} (matrix
 331 of $K \times d$ with each chain state as row vector), using differential evolution [*Storn and Price,*
 332 *1997; Price et al., 2005*]. If A is a subset of δ -dimensional space of the original parameter
 333 space, $\mathbb{R}^\delta \subseteq \mathbb{R}^d$, then a jump in the k th chain, $k = \{1, \dots, K\}$ at iteration $t = \{2, \dots, T\}$
 334 is calculated using

$$\begin{aligned} \Delta_{k,A}^* &= (\mathbf{1}_\delta + \boldsymbol{\lambda})\gamma(\delta) \left[\sum_{j=1}^{\tau} \boldsymbol{\theta}_{\mathbf{g}_j,A}^{t-1} - \sum_{j=1}^{\tau} \boldsymbol{\theta}_{\mathbf{r}_j,A}^{t-1} \right] + \boldsymbol{\zeta} \\ \Delta_{k,\neq A}^* &= 0, \end{aligned} \tag{14}$$

335 where $\gamma = 2.38/\sqrt{2\tau D}$ is the jump rate, τ denotes the number of chain pairs used to gen-
 336 erate the jump, and \mathbf{g} and \mathbf{r} are τ -vectors with integer values drawn without replacement
 337 from $\{1, \dots, k-1, k+1, \dots, K\}$. The values of $\boldsymbol{\lambda}$ and $\boldsymbol{\zeta}$ are sampled independently from
 338 $\mathcal{U}_\delta(-c, c)$ and $\mathcal{N}_\delta(0, c^*)$ with, typically, $c = 0.1$ and c^* small compared to the width of the
 339 target distribution, $c^* = 10^{-12}$ say.

340 The candidate point of chain k at iteration t then becomes

$$\boldsymbol{\theta}_k^* = \boldsymbol{\theta}_k^{t-1} + \Delta_k^*, \tag{15}$$

341 and the Metropolis ratio is used to determine whether to accept this proposal or not. The
 342 DREAM algorithm solves an important practical problem in MCMC simulation, namely
 343 that of choosing an appropriate scale and orientation of the proposal distribution. Section
 344 3.3 will detail the procedure for selecting the subset A of dimensions of the parameter
 345 space that will be updated each time a proposal is created.

We now proceed with a pseudo-code of DREAM_(ABC).

Algorithm 3 DREAM_(ABC)-Markov chain Monte Carlo sampler

- 1: At iteration $t = 1$,
- 2: **for** $k = 1, \dots, K$ **do**
- 3: Sample $\boldsymbol{\theta}_k^1$ from the prior, $\boldsymbol{\theta}_k^1 \sim p(\boldsymbol{\theta})$
- 4: Simulate data \mathbf{Y} using $\boldsymbol{\theta}_k^1$, $\mathbf{Y} \sim \text{Model}(\boldsymbol{\theta}_k^1)$
- 5: Calculate the fitness, $f(\boldsymbol{\theta}_k^1) = \epsilon - \rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}_k^1)))$
- 6: **end for**
- 7: At iteration $t > 1$,
- 8: **for** $t = 2, \dots, T$ **do**
- 9: **for** $k = 1, \dots, K$ **do**
- 10: Determine subset A , the dimensions of the parameter space to be updated.
- 11: Calculate the jump vector, $\boldsymbol{\Delta}_k^*$ using Equation (14)
- 12: Compute the proposal, $\boldsymbol{\theta}_k^* = \boldsymbol{\theta}_k^{t-1} + \boldsymbol{\Delta}_k^*$
- 13: Simulate data \mathbf{Y} using $\boldsymbol{\theta}_k^*$, $\mathbf{Y} \sim \text{Model}(\boldsymbol{\theta}_k^*)$
- 14: Calculate the fitness, $f(\boldsymbol{\theta}_k^*) = \epsilon - \rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}_k^*)))$
- 15: Calculate the acceptance probability using Equation (13),

$$p_{\text{acc}}(\boldsymbol{\theta}_k^*) = \max\left(\mathbb{I}(f(\boldsymbol{\theta}_k^*) \geq f(\boldsymbol{\theta}_k^{t-1})), \mathbb{I}(f(\boldsymbol{\theta}_k^*) \geq 0)\right)$$

- 16: If $p_{\text{acc}}(\boldsymbol{\theta}_k^*) = 1$, set $\boldsymbol{\theta}_k^t = \boldsymbol{\theta}_k^*$ otherwise remain at current state, $\boldsymbol{\theta}_k^t = \boldsymbol{\theta}_k^{t-1}$
 - 17: **end for**
 - 18: Compute \hat{R} -statistic for each entry of $\boldsymbol{\theta}$ using last 50% of samples in each chain.
 - 19: If $\hat{R}_j \leq 1.2$ for all $j = \{1, \dots, d\}$ stop, otherwise continue.
 - 20: **end for**
-

346

347 In summary, DREAM_(ABC), runs K different Markov chains in parallel. Multivariate
 348 proposals in each chain $k = \{1, \dots, K\}$ are generated by taking a fixed multiple of the

349 difference of two or more randomly chosen members (chains) of Θ (without replacement).
 350 By accepting each jump with binary probability of Equation (13) a Markov chain is
 351 obtained, the stationary or limiting distribution of which is the posterior distribution.
 352 Because the joint pdf of the K chains factorizes to $\pi(\boldsymbol{\theta}_1) \times \dots \times \pi(\boldsymbol{\theta}_K)$, the states of
 353 the individual chains are independent at any iteration after DREAM_(ABC) has become
 354 independent of its initial value. After this burn-in period, the convergence of DREAM_(ABC)
 355 can thus be monitored with the \hat{R} -statistic of *Gelman and Rubin* [1992].

356 The jump distribution in Equation (14) of DREAM_(ABC) is easily implemented in ABC-
 357 PMC to help generate trial moves. This could further improve the scale and orientation of
 358 the proposals, but comes at an increased computational cost. The conditional probability
 359 to move from $\boldsymbol{\theta}_{j-1}^u$ to $\boldsymbol{\theta}_j^i$ or $q_d(\boldsymbol{\theta}_{j-1}^u \rightarrow \boldsymbol{\theta}_j^i)$ in Equation (4) of ABC-PMC is easy to calculate
 360 for a (multi)normal proposal distribution, but requires significantly more CPU-resources
 361 if the jumping kernel of DREAM_(ABC) is used. Let's assume, for instance, that $\tau = 1$ and
 362 $\boldsymbol{\lambda} = \mathbf{0}_d$ in Equation (14). The probability to transition in ABC-PMC from the current
 363 state, $\boldsymbol{\theta}_{j-1}^u$, to the proposal, $\boldsymbol{\theta}_j^i$, is then equivalent to

$$q(\boldsymbol{\theta}_j^i | \boldsymbol{\theta}_{j-1}^u) = \sum_{m=1}^N \sum_{o=1}^N \psi(\boldsymbol{\theta}_{j-1}^u + \gamma(\boldsymbol{\theta}_{j-1}^m - \boldsymbol{\theta}_{j-1}^o) | \boldsymbol{\Sigma}_d); i \neq o \neq m \quad (16)$$

364 where ψ denotes the normal density with covariance matrix $\boldsymbol{\Sigma}_d = (c^*)^2 I_d$. This Equation
 365 is of computational complexity $\mathcal{O}(N^2)$ and becomes particularly CPU-intensive for large
 366 N and/or if more than one pair of chains ($\tau > 1$) is used to create proposals. More
 367 fundamentally, the lack of subspace sampling (see next section) in ABC-PMC deteriorates
 368 search efficiency in high-dimensional spaces (shown later). We therefore do not consider
 369 this alternative jumping distribution in ABC-PMC.

3.3. Randomized subspace sampling

Subspace sampling is implemented in DREAM_(ABC) by only updating randomly selected dimensions of θ_k^{t-1} each time a proposal is generated. Following the default of the DREAM suite of algorithms [Vrugt *et al.*, 2008, 2009, 2011; Laloy and Vrugt, 2012] we use a geometric series of n_{CR} different crossover values and store this in a vector, $\text{CR} = \{\frac{1}{n_{\text{CR}}}, \frac{2}{n_{\text{CR}}}, \dots, 1\}$. The prior probability of each crossover value is assumed equal and defines a vector \mathbf{p} with n_{CR} copies of $\frac{1}{n_{\text{CR}}}$. We create the set A with selected dimensions to be updated as follows

Algorithm Subspace sampling

```

1: for  $k = 1, \dots, K$  do
2:   Define  $A$  to be an empty set,  $A = \emptyset$ 
3:   Sample  $P$  from the discrete multinomial distribution,  $P \sim \mathfrak{F}(\text{CR}|\mathbf{p})$ 
4:   Draw  $d$  labels from a multivariate uniform distribution,  $\mathbf{Z} \sim \mathcal{U}_d[0, 1]$ 
5:   for  $j = 1, \dots, d$  do
6:     if  $\mathbf{Z}_j > (1 - P)$  then
7:       Add dimension  $j$  to  $A$ 
8:     end if
9:   end for
10:  if  $A = \emptyset$  then
11:    Choose one index from  $\{1, \dots, d\}$  and add to  $A$ 
12:  end if
13: end for

```

The number of dimensions stored in A ranges between 1 and d and depends on the sampled value of the crossover. This relatively simple randomized selection strategy enables single-site Metropolis sampling (one dimension at a time), Metropolis-within-Gibbs (one or a group of dimensions) and regular Metropolis sampling (all dimensions). To enhance search efficiency, the probability of each n_{CR} crossover values is tuned adaptively during burn-in by maximizing the normalized Euclidean distance between successive states of the

383 K chains [Vrugt *et al.*, 2009]. The only algorithmic parameter that needs to be defined
384 by the user is n_{CR} , the number of crossover values used. We use the standard settings
385 of DREAM and use $n_{\text{CR}} = 3$ in all the calculations reported herein. This concludes the
386 algorithmic description of DREAM_(ABC).

4. Numerical experiments

387 The next section compares the efficiency of ABC-REJ, ABC-PMC and DREAM_(ABC)
388 for two synthetic and one real-world experiment. These case studies cover a diverse
389 set of problem features, including high-dimensionality, nonlinearity, non-convexity, and
390 numerous local optima. Perhaps not surprisingly, our trials with ABC-REJ show that
391 rejection sampling is highly inefficient when confronted with multidimensional parameter
392 spaces, unless the prior sampling distribution closely mimics the target distribution of
393 interest. In practice, this is an unreasonable expectation and we therefore discard the
394 ABC-REJ algorithm after the first case study and focus our attention on the results of
395 ABC-PMC and DREAM_(ABC).

396 In all our calculations with ABC-PMC we create $N = 1,000$ samples at each iteration,
397 $j = \{1, \dots, J\}$ using values of ϵ that are listed in each case study and have been determined
398 through trial-and-error (e.g. *Sadegh and Vrugt [2013]*). In DREAM_(ABC) we need to define
399 the number of chains, K and the total number of function evaluations, $M = K \cdot T$. Their
400 values are listed in each individual case study. For all the other algorithmic variables in
401 DREAM_(ABC) we use standard settings recommended in *Vrugt et al. [2009]*.

4.1. Synthetic benchmark experiments: Gaussian mixture model

402 We now benchmark the performance of DREAM_(ABC) by application to the Gaussian
 403 mixture model in Equation (5). We set $\epsilon = 0.025$ and run $K = 10$ different chains using
 404 a total of $M = 50,000$ function evaluations. Fig. 2c displays the marginal distribution of
 405 the posterior samples using a burn-in of 50%. It is evident that the adaptive capabilities
 406 of DREAM_(ABC) enables it to track the target distribution. The density of samples in the
 407 tails has somewhat improved considerably compared to ABC-PMC. The burn-in required
 408 for multi-chain methods such as DREAM_(ABC) is relatively costly for this one-dimensional
 409 problem, and hence it is not surprising that rejection sampling provides a nicer sample
 410 of the mixture distribution. The DREAM_(ABC) approximation of the target is easily
 411 enhanced by creating more samples.

412 Table 1 lists the acceptance rate (AR, %), number of function evaluations, and ϵ value
 413 used with ABC-REJ, ABC-PMC and DREAM_(ABC). It is evident that DREAM_(ABC) is
 414 most efficient in sampling the target distribution. The acceptance rate of DREAM_(ABC)
 415 is between 3 to 6 times higher than that ABC-PMC and ABC-REJ, respectively.

4.2. Synthetic benchmark experiments: 20-dimensional bivariate distribution

416 A more demanding test of the ABC-PMC and DREAM_(ABC) algorithms can be devised
 417 by using a multi-dimensional target distribution. We consider a set of ten bivariate normal
 418 distributions

$$h_i(\mu_{i,1}, \mu_{i,2}) \sim \mathcal{N}_2 \left(\begin{bmatrix} \mu_{i,1} \\ \mu_{i,2} \end{bmatrix}, \begin{bmatrix} 0.01^2 & 0 \\ 0 & 0.01^2 \end{bmatrix} \right), \quad (17)$$

419 with unknown mean of the i th component, $\boldsymbol{\mu}_i = \{\mu_{i,1}, \mu_{i,2}\}$ and fixed covariance matrix.
 420 We now generate $n = 20$ observations by sampling the mean of each bivariate distribution
 421 from $\mathcal{U}_2[0, 10]$. The "observed" data are plotted in Figure 3 using the red cross symbols.

Each of the bivariate means is now subject to inference with ABC, which results in a $d = 20$ dimensional parameter estimation problem. The simulated data, \mathbf{Y} are created by evaluating Equation (17) fifty different times for each proposal, $\boldsymbol{\theta}^* = [\boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_{10}^*]$. The mean of the fifty samples from each bivariate distribution is stored in \mathbf{Y} (10 by 2 matrix) and compared to the observed data using the following distance function [Turner and Sederberg, 2012]

$$\rho(\tilde{\mathbf{Y}}, \mathbf{Y}(\boldsymbol{\theta}^*)) = \sqrt{\frac{1}{20} \sum_{i=1}^{10} \sum_{j=1}^2 (\tilde{\mathbf{Y}}_{(i,j)} - \mathbf{Y}_{(i,j)}(\boldsymbol{\theta}^*))^2} \quad (18)$$

We assume a noninformative (uniform) prior, $\boldsymbol{\theta} \sim \mathcal{U}_{20}[0, 10]$ and set $\boldsymbol{\epsilon} = \{3, 2.5, 2.1, 1.8, 1.6, 1.3, 1.1, 0.9, 0.8, 0.7, 0.6\}$ (ABC-PMC) and $\epsilon = 0.025$, $K = 15$, and $M = 200,000$ (DREAM_(ABC)). The results of the analysis are presented in Table 2 and Figure 3.

Figure 3 plots the posterior samples (plus symbol) derived from (a) ABC-PMC and (b) DREAM_(ABC). The sampled solutions cluster around the observed means (red cross) of the bivariate normal distributions. The size of the posterior uncertainty differ markedly between both algorithms. The posterior samples of DREAM_(ABC) are in excellent agreement with the bivariate target distributions. The samples group tightly around their observed counterparts, and their structure is in excellent agreement with the covariance matrix of the target distribution. The ABC-PMC samples, on the contrary, exhibit too much scatter. This finding is not surprising. The ABC-PMC sampler terminated its search prematurely with values of $\rho(\cdot)$ between 0.5 and 0.6. This threshold is much larger than the value of $\epsilon = 0.025$ required to converge to the target distribution. Subspace sampling

442 would significantly enhance the results of ABC-PMC (not shown), but such modifications
443 could affect the theoretical convergence properties.

444 The DREAM_(ABC) algorithm not only better recovers the actual target distribution, but
445 its sampling efficiency is also superior. To illustrate this in more detail, consider Table
446 2 that lists the acceptance rate (AR, %), number of function evaluations, and ϵ value of
447 ABC-PMC and DREAM_(ABC). The acceptance rate of DREAM_(ABC) of about 19.72% is
448 more than 1,000 times higher than that of ABC-PMC (0.019%). This marks a three-order
449 of magnitude improvement in search efficiency, which in large part is due to the ability
450 of DREAM_(ABC) to sample one or groups of variables in turn. This conditional sampling
451 is necessary to traverse multi-dimensional parameter spaces in pursuit of the posterior
452 distribution.

453 The present case study clearly illustrates the advantages of DREAM_(ABC) when con-
454 fronted with multi-dimensional parameter spaces. The algorithm requires about $M =$
455 200,000 function evaluations to successfully recover the 20-d target distribution. To illus-
456 trate this in more detail, consider Figure 4 which displays trace plots of the \hat{R} -statistic
457 of *Gelman and Rubin* [1992] using the last 50% of the samples stored in each of the K
458 chains. This convergence diagnostic compares for each parameter the between- and within-
459 variance of the chains. Because of asymptotic independence, the between-member variance
460 and \hat{R} -diagnostic can be estimated consistently from a single DREAM_(ABC) trial. Values
461 of \hat{R} smaller than 1.2 indicate convergence to a limiting distribution. The DREAM_(ABC)
462 algorithm needs about 40,000 function evaluations to officially reach convergence and
463 generate a sufficient sample of the posterior distribution. This is much less than the
464 $M = 200,000$ function evaluations used in this study. Obviously, one should be careful

465 to judge convergence of the sampled Markov chains based on a single diagnostic, yet vi-
 466 sual inspection of the sampled trajectories confirms an adequate mixing of the different
 467 chains and converge of DREAM_(ABC) after about 15,000 function evaluations. This is sub-
 468 stantially lower than the approximately 40,000 function evaluations estimated with the
 469 \hat{R} -statistic, simply because the second half of the chain is used for monitoring convergence.

4.3. Hydrologic modeling: Sacramento Soil Moisture Accounting Model

470 A more realistic case study is now devised and used to illustrate the advantages
 471 DREAM_(ABC) can offer in real-world modeling problems. We consider simulation of the
 472 rainfall-runoff transformation using the SAC-SMA conceptual hydrologic model. This
 473 model has been developed by *Burnash et al.* [1973] and is used extensively by the Na-
 474 tional Weather Service for flood forecasting throughout the United States. The model
 475 has been described in detail in many previous publications, and we therefore summarize
 476 in Table 3 the fourteen parameters that require calibration and their prior uncertainty
 477 ranges.

478 Daily data of mean areal precipitation, $\tilde{\mathbf{P}} = \{\tilde{p}_1, \dots, \tilde{p}_n\}$, mean areal potential evapo-
 479 ration and streamflow, $\tilde{\mathbf{Y}} = \{\tilde{y}_1, \dots, \tilde{y}_n\}$, from the French Broad River basin at Asheville,
 480 North Carolina are used in the present study. In keeping with *Vrugt and Sadegh* [2013] we
 481 use the annual runoff coefficient, $S_1(\tilde{\mathbf{Y}})$ (-), the annual baseflow index, $S_2(\tilde{\mathbf{Y}})$ (-), and the
 482 flow duration curve, $S_3(\tilde{\mathbf{Y}})$ (day/mm) and $S_4(\tilde{\mathbf{Y}})$ (-), as summary metrics of the discharge
 483 data. A detailed description of each summary statistic is given by *Vrugt and Sadegh* [2013]
 484 and interested readers are referred to this publication for further details. This leaves us
 485 with $L = 4$ four summary statistics for three different hydrologic signatures.

486 We use the following composite distance function to quantify the distance between the
 487 observed and simulated summary statistics

$$\rho(S(\tilde{\mathbf{Y}}), S(\mathbf{Y}(\boldsymbol{\theta}^*))) = \max(|S_i(\tilde{\mathbf{Y}}) - S_i(\mathbf{Y}(\boldsymbol{\theta}^*))|) \quad i = \{1, \dots, L\}, \quad (19)$$

488 and help determine whether to accept $\boldsymbol{\theta}^*$ or not. Model simulations that simultaneously
 489 satisfy each of the four summary metrics within their tolerance thresholds are considered
 490 behavioral, and hence constitute samples from the posterior distribution. This compos-
 491 ite formulation differs fundamentally from multi-criteria model calibration approaches in
 492 which the summary statistics (objective functions) are assumed non-commensurate and
 493 therefore treated independently from one another. This latter approach gives rise to a
 494 Pareto solution set (rather than posterior distribution) and quantifies the trade-offs in the
 495 fitting of the different metrics.

496 To maximize the search efficiency of ABC-PMC we use $\epsilon = \{1, 0.3, 0.15, 0.1, 0.07, 0.06, 0.04, 0.025\}$.
 497 This sequence is determined from a preliminary run of ABC-PMC with adaptive selec-
 498 tion of ϵ (see Appendix B of *Sadegh and Vrugt* [2013]). The DREAM_(ABC) sampler is
 499 executed using default values of the algorithmic parameters and $\epsilon = 0.025$, $K = 15$, and
 500 $M = 200,000$. Tables 3 and 4 and Figures 5, 6, and 7 summarize our main findings.

501 Table 4 compares the computational efficiency of ABC-PMC and DREAM_(ABC). For
 502 completeness, we also list the results of DREAM using a residual-based Gaussian like-
 503 lihood function, hereafter referred to as DREAM_(RBGL). The DREAM_(ABC) algorithm
 504 has an acceptance rate (AR, %) of about 3.14% and requires 200,000 SAC-SMA model
 505 evaluations to generate 40,000 posterior samples. The ABC-PMC sampler, on the con-
 506 trary, is far less efficient (AR = 0.046%) and needs about 2.2 million function evaluations

507 to produce 1,000 posterior samples. This constitutes a more than ten times difference in
508 sampling efficiency, and favors the use of DREAM_(ABC) for diagnostic inference of complex
509 and CPU-intensive models.

510 The acceptance rate of DREAM_(RBGL) of 4.54% is much lower than the theoretical
511 (optimal) value of about 23.4% for the considered dimensionality of the target distribution.
512 This finding is not surprising and can be explained by the non-ideal properties of the
513 SAC-SMA response surface [*Duan et al.*, 1992], which, to a large extent, are inflicted
514 by poor numerics [*Clark and Kavetski*, 2010; *Kavetski and Clark*, 2010; *Schoups and*
515 *Vrugt*, 2010]. The use of an explicit, Euler-based, integration method introduces pits and
516 local optima (amongst others) on the response surface, and their presence deteriorates
517 the search efficiency of MCMC methods. An implicit, time-variable, integration method
518 would give a smoother response surface but at the expense of an increase in CPU-time.
519 This increase in computational cost, will however, be balanced by a decrease in the number
520 of model evaluations needed for a MCMC algorithm to converge to a limiting distribution.

521 Figure 5 presents histograms of the marginal posterior distributions derived with ABC-
522 PMC (top panel), DREAM_(ABC) (middle panel) and DREAM_(RBGL) (lower panel). We
523 display the results of a representative set of six SAC-SMA parameters and plot, from
524 left to right across each panel, the posterior distributions of PCTIM, ADIMP, LZFSM,
525 LZFPM, and LZPK. The x -axis matches exactly the ranges of each parameter used in the
526 (uniform) prior distribution.

527 The marginal distributions derived from both sampling methods are in good agree-
528 ment, and exhibit similar functional shapes. This inspires confidence in the ability of
529 DREAM_(ABC) to correctly sample the target distribution. Most histograms extent a large

530 part of the prior distribution, which suggests that the parameters are poorly defined by
531 calibration against the four different summary statistics. This finding is perhaps not sur-
532 prising. The four metrics used in this study are not sufficient, and extract only a portion
533 of the information available in the discharge calibration data set. We will revisit this
534 issue in the final paragraph of this section. Information theory can help to determine an
535 approximate set of sufficient statistics, but this is beyond the scope of the present paper.

536 We can further constrain the behavioral (posterior) parameter space by adding other
537 signatures of catchment behavior to the current set of summary metrics. But, it is not
538 particularly clear whether this would actually support the purpose of diagnostic model
539 evaluation in which the (our) goal is not to just find the best possible fit of some model to
540 some data set, but rather to detect and pinpoint (epistemic) errors arising from inadequate
541 or incomplete process representation. The chosen metrics appear relatively insensitive to
542 rainfall data errors (not shown herein) and therefore exhibit useful diagnostic power. The
543 bottom panel illuminates what happens to the SAC-SMA parameters if a least-squares
544 type likelihood function is used for posterior inference. The parameters appear to be
545 much better resolved by calibration against the observed discharge data but the remain-
546 ing error residuals violate assumptions of homoscedasticity, normality and independence
547 (not shown in detail). In part, this is due to a lack of treatment of rainfall data errors,
548 whose probabilistic properties are difficult to accurately represent in a likelihood function.
549 The generalized likelihood function of *Schoups and Vrugt* [2010] provides ways to handle
550 nontraditional residual distributions, nevertheless, this approach does not separate the
551 contribution of individual error sources, and is therefore unable to provide insights into
552 model malfunctioning. Note that the histograms of PCTIM, ADIMP, LZFSM, LZFPM,

553 and LZPK are relatively tight and well described by a normal distribution, except for
554 PCTIM which is hitting its lower bound.

555 To illustrate how the SAC-SMA posterior parameter uncertainty translates into modeled
556 discharge uncertainty, please consider Figure 6 that presents time series plots of the 95%
557 streamflow simulation uncertainty ranges (gray region) for a selected 500-day portion
558 of the 3-year evaluation period derived from the posterior samples of ABC-PMC (top
559 panel) and DREAM_(ABC) (middle panel). The observed discharge data are indicated with
560 solid circles. Both time-series plots are in excellent agreement with simulated discharge
561 dynamics that appear visually very similar and uncertainty ranges that envelop a large
562 majority of the streamflow observations. Epistemic errors are not readily visible, yet
563 this requires much further analysis possibly with the use of additional summary metrics.
564 Previous results for this data set presented in *Vrugt and Sadegh* [2013], demonstrated an
565 inability of the 7-parameter hmodel [*Schoups and Vrugt*, 2010] to simulate accurately the
566 immediate response of the watershed to rainfall. This structural error can be resolved
567 by model correction, a topic that will be studied in future publications. Note that the
568 posterior mean RMSE derived with DREAM_(ABC) (0.72 mm/day) is somewhat lower than
569 its counterpart from ABC-PMC (0.81 mm/day). This inconsistency conveys a difference
570 in sampling density and posterior approximation.

571 For completeness, the bottom panel of Figure 6 plots the 95% streamflow simulation
572 uncertainty ranges derived from formal Bayes using a least-squares likelihood function. To
573 enable a direct comparison with the results for diagnostic inference in the top two panels,
574 we only consider the effect of parameter uncertainty on simulated discharge dynamics.
575 The coverage has decreased substantially to about 13%, which is hardly surprising given

576 the relatively small width of the marginal distributions shown in Fig. 5. The RMSE of
577 the posterior mean SAC-SMA simulation (0.55 mm/day) is considerably lower than its
578 counterparts derived from ABC-PMC and DREAM_(ABC). This finding is not alarming but
579 warrants some discussion. The four summary metrics used for diagnostic inference only
580 extract partial information from the available discharge observations. This insufficiency
581 makes it difficult to find a posterior model that "best" fits, in least-squares sense, the
582 streamflow data, which is expected from a Gaussian likelihood function with homoscedas-
583 tic measurement error. Also, the main purpose of diagnostic model evaluation with ABC
584 is not that of model calibration, but rather to provide insights into model malfunctioning.
585 Residual-based model calibration approaches provide little guidance on this issue which
586 limits our ability to learn from the calibration data.

587 Table 5 presents summary variables (coverage, width, root mean square error, bias,
588 and correlation coefficient) of the performance of the posterior mean SAC-SMA discharge
589 simulation derived from the samples of ABC-PMC, DREAM_(ABC) and DREAM_(RBGL).
590 We list results for the 5-year calibration and 3-year evaluation period. These statistics
591 confirm our previous findings. The two different ABC sampling methods provide very
592 similar results, and exhibit a better coverage of the discharge observations, larger width
593 of the 95% simulation uncertainty ranges, and higher posterior mean RMSE than least-
594 squares inference. The performance of the SAC-SMA model does not deteriorate during
595 the evaluation period. In fact, the RMSE of the ABC derived posterior mean simulation
596 substantially improves during the evaluation period, whereas this is not the case with
597 least-squares fitting. This is a heartening prospect, and suggests (among others) that the

598 chosen summary metrics at least partially represent the underlying signatures of watershed
599 behavior.

600 To provide more insights into the convergence behavior of DREAM_(ABC), Figure 7 plots
601 the evolution of the \hat{R} -statistic of *Gelman and Rubin* [1992]. Each of the SAC-SMA
602 parameters is coded with a different color. About 160,000 SAC-SMA model evaluations
603 are required to converge to a limiting distribution. This marks a significant improvement in
604 sampling efficiency over the ABC-PMC sampler which requires about 2.2 million function
605 evaluations to create 1,000 posterior samples.

606 We now turn our attention to the simulated values of the summary metrics. Figure 8
607 plots histograms of the posterior summary statistics derived with ABC-PMC (top panel)
608 and DREAM_(ABC) (middle panel). The observed values of summary statistics are sep-
609 arately indicated in each plot with a red cross. The marginal distributions of summary
610 metrics generally center around their measured values with the exception of the histogram
611 of S_1 (annual runoff coefficient) that appears heavily skewed to the right. This points to a
612 potential deficiency in the SAC-SMA model structure, yet this requires further analysis.
613 For completeness, the bottom panel plots the posterior summary metric distributions de-
614 rived from a residual-based likelihood function. This approach provides the closest fit to
615 the observed streamflow data, but at the expense of summary metrics S_3 and S_4 (flow du-
616 ration curve) that deviate considerably from their observed values. This finding highlights
617 a profound difference in methodology between diagnostic model evaluation with ABC and
618 residual-based inference. Of course, we could have used a different, and perhaps more
619 reasonable, likelihood function for the error residuals [e.g. *Schoups and Vrugt* [2010]].

620 This would affect some of our findings in Figure 8. Nevertheless, this is outside the scope
621 of the present paper, and we leave such comparison for future work.

622 Finally, Figure 9 presents bivariate scatter plots of the posterior samples generated with
623 DREAM_(ABC). We display the results for a representative set of all SAC-SMA parameter
624 pairs including (A) PCTIM - LZPK, (B) LZSK - LZPK, (C) PCTIM - LZFP, (D)
625 PCTIM - LZSK, (E) LZPK - PFREE, and (F) ADIMP - LZPK. The axes in each of
626 the six plots are in agreement with those used in the prior distribution. The bivariate
627 posterior samples are confined to a densely-sampled rectangular (or square) space, and
628 occupy a significant portion of the prior distribution. The white area immediately outside
629 of the sampled space is made up of nonbehavioral solutions with fitness values smaller
630 than zero (at least one summary metric is ϵ removed from its measured counterpart). The
631 binary acceptance rule used in DREAM_(ABC) introduces a rather sharp demarcation of the
632 behavioral solution space, nevertheless the sampled posterior distribution is in excellent
633 agreement with its counterpart derived from ABC-PMC (see Figure 5). One could argue
634 that that for this type of target distribution uniform random sampling should suffice.
635 However, this method, which is at the heart of ABC-REJ, is highly inefficient in multi-
636 dimensional parameter spaces. Many millions of function evaluations would be needed
637 to provide a sufficient sample of the posterior distribution. This is rather cumbersome,
638 particularly if, as in the present case (not shown), the posterior samples exhibit parameter
639 correlation.

640 The focus of the present paper has been on improving ABC sampling efficiency to per-
641 mit diagnostic inference of complex system models involving multi-dimensional parameter
642 and summary metric spaces. Subsequent work can now focus on the intended purpose of

643 diagnostic model evaluation and that is to help detect, diagnose, and resolve model struc-
644 tural deficiencies [*Vrugt and Sadegh, 2013*]. Commonly used model-data fusion approaches
645 provide limited guidance on this important issue, in large part because of their aggregated
646 treatment of input (forcing) data and epistemic errors. The use of summary statistics for
647 statistical inference holds great promise, not only because signatures of system behav-
648 ior are much less sensitive to, for instance, precipitation data errors than residual-based
649 model fitting approaches, but also because the metrics can be devised in such a way that
650 they relate directly to individual process descriptions and thus model components. This
651 has important diagnostic advantages. Failure to fit one or more summary metrics can be
652 directly addressed through correction of the responsible model component(s). This itera-
653 tive process of inference, adjustment and refinement constitutes the basis of the scientific
654 method. This also leaves the possibility to collect additional data to help validate new
655 components of the model.

656 A recurrent issue with the application of diagnostic inference will be sufficiency of the
657 summary metrics. Ideally, the summary metrics contain as much information as the
658 original data itself. Unfortunately, for most systems it will be rather difficult to find
659 a set of sufficient summary statistics, unless each calibration data measurement is used
660 as independent metric (e.g. *Sadegh and Vrugt [2013]*) but this defeats the purpose of
661 diagnostic inference. Actually, it is not particularly clear whether sufficiency of the metrics
662 is required to help detect and resolve epistemic errors. If deemed necessary, then one
663 possible solution is to adapt formal Bayes and to use the summary metrics as an explicit
664 prior. This type of approach has shown to significantly enhance the results of geophysical
665 inversion [*Lochbühler et al., 2014*].

5. Summary and Conclusions

666 The paper by *Vrugt and Sadegh* [2013] has introduced approximate Bayesian compu-
667 tation (ABC) as vehicle for diagnostic model evaluation. Successful application of this
668 methodology requires availability of an efficient sampling method that rapidly explores
669 the space of behavioral models. Commonly used rejection sampling approaches adopt
670 a boxcar kernel (0/1) to differentiate between behavioral ("1") and non-behavioral ("0")
671 solutions, and use full-dimensional updating in pursuit of the posterior parameter distri-
672 bution. This approach might work well for low-dimensional problems (e.g. $d \leq 10$) but
673 is not particularly efficient in high-dimensional parameter spaces which require partial
674 (subspace) sampling to rapidly locate posterior solutions.

675 In this paper, we have introduced DREAM_(ABC) to permit diagnostic inference of
676 complex system models. This approach uses Metropolis-within-Gibbs simulation with
677 DREAM [*Vrugt et al.*, 2008, 2009] to delineate the space of behavioral (posterior) mod-
678 els. Three different case studies involving a simple one-dimensional toy problem, a 20-
679 dimensional mixture of bivariate distributions, and a 14-dimensional hydrologic model
680 calibration problem illustrate that DREAM_(ABC) is about 3 - 1,000 times more efficient
681 than commonly used ABC sampling approaches. This gain in sampling efficiency increases
682 with dimensionality of the parameter space.

683 The source code of DREAM_(ABC) is written in MATLAB and available upon request
684 from the second author: `jasper@uci.edu`. This code includes (amongst others) the three
685 different case studies considered herein and implements many different functionalities
686 (postprocessing and visualization tools, convergence and residual diagnostics) to help users
687 analyze their results.

688 **Acknowledgments.** Both authors highly appreciate the support and funding from the
689 UC-Lab Fees Research Program Award 237285. The comments of the three anonymous
690 referees have improved the current version of this manuscript.

References

- 691 Barnes, C., S. Filippi, M.P.H. Stumpf, and T. Thorne (2011), Con-
692 siderate approaches to achieving sufficiency for ABC model selection,
693 <http://arxiv.org/pdf/1106.6281v2.pdf>, 1–21.
- 694 Beaumont, M.A., W. Zhang, and D.J. Balding (2002), Approximate Bayesian computation
695 in population genetics, *Genetics*, *162*(4), 2025–2035.
- 696 Beaumont, M.A., J.M. Cornuet, J.M., Marin, and C.P. Robert (2009), Adaptive approx-
697 imate Bayesian computation, *Biometrika*, *asp052*, 1–8.
- 698 Beaumont, M.A. (2010), Approximate Bayesian Computation in Evolution and Ecology,
699 *Annual Review of Ecology, Evolution, and Systematics*, *41*, 379–406.
- 700 Bertorelle, G., A. Benazzo, and S. Mona (2010), ABC as a flexible framework to estimate
701 demography over space and time: some cons, many pros, *Molecular Ecology*, *19*, 2609–
702 2625.
- 703 Blum, M.G.B., and O. François (2010), Non-linear regression models for approximate
704 Bayesian computation, *Statistics and Computing*, *20*, 63–73.
- 705 Burnash, R.J., R.L. Ferral, R.L., and R.A. McGuire (1973), A generalized streamflow sim-
706 ulation system: Conceptual modeling for digital computers, Joint Federal-State River
707 Forecast Center, Sacramento, CA, USA.
- 708 Clark, M.P., and D. Kavetski (2010), Ancient numerical demons of conceptual hydro-
709 logical modeling: 1. Fidelity and efficiency of time stepping schemes, *Water Resources*
710 *Research*, *46*, W10510, doi:10.1029/2009WR008894.
- 711 Csilléry, K., M.G.B. Blum, O.E. Gaggiotti, and O. François (2010), Approximate Bayesian
712 Computation (ABC) in practice, *Trends in Ecology & Evolution*, *25*, 410–418.

- 713 Del Moral, P., A. Doucet, and A. Jasra (2008), An adaptive sequential Monte Carlo
714 method for approximate Bayesian computation, Technical Report, Imperial College
715 London.
- 716 Diggle, P.J., and R.J. Gratton (1984), Monte Carlo methods of inference for implicit
717 statistical models, *Journal of the Royal Statistical Society Series B*, *46*, 193–227.
- 718 Duan, Q., S. Sorooshian, and V. Gupta (1992), Effective and efficient global optimization
719 for conceptual rainfall-runoff models, *Water Resources Research*, *28*(4), 1015–1031.
- 720 Duan, Q., J. J. Schaake, V. Andréassian, S. Franks, G. Goteti, H.V. Gupta, Y.M. Gu-
721 sev, F. Habets, A. Hall, L. Hay, T. Hogue, M. Huang, G. Leavesley, X. Liang, O.N.
722 Nasonova, J. Noilhan, L. Oudin, S. Sorooshian, T. Wagener, and E.F. Wood (2006),
723 Model parameter estimation experiment (MOPEX): An overview of science strategy
724 and major results from the second and third workshops, *Journal of Hydrology*, *320*,
725 3–17.
- 726 Evin, G., D. Kavetski, M. Thyer, and G. Kuczera (2013), Pitfalls and improvements
727 in the joint inference of heteroscedasticity and autocorrelation in hydrological model
728 calibration, *Water Resources Research*, *49*, doi:10.1002/wrcr.20284.
- 729 Gelman, A., D.B. Rubin (1992), Inference from iterative simulation using multiple se-
730 quences, *Statistical Science*, *7*, 457–472.
- 731 Grelaud, A., C. Robert, J. Marin, F. Rodolphe, and J. Taly (2009) ABC likelihood-free
732 methods for model choice in Gibbs random fields, *Bayesian Analysis*, *4* (2), 317–336.
- 733 Gupta, H.V., T. Wagener, and Y. Liu (2008), Reconciling theory with observations: el-
734 ements of a diagnostic approach to model evaluation, *Hydrological Processes*, *22* (18),
735 3802–3813.

- 736 Haario, H., E. Saksman, and J. Tamminen (1999), Adaptive proposal distribution for
737 random walk Metropolis algorithm, *Comp. Stat.*, *14*(3), 375–395.
- 738 Haario, H., E. Saksman, and J. Tamminen (2001), An adaptive Metropolis algorithm,
739 *Bernoulli*, *7*, 223–242.
- 740 Hastings, H. (1970), Monte Carlo sampling methods using Markov chains and their ap-
741 plications, *Biometrika*, *57*, 97–109.
- 742 Higuchi, T. (1997), Monte Carlo filter using the genetic algorithm operators, *Journal of*
743 *Statistical Computation and Simulation*, *59*, 1–23.
- 744 Joyce, P. and P. Marjoram (2008), Approximately sufficient statistics and Bayesian com-
745 putation, *Statistical Applications in Genetics and Molecular Biology*, *7* (1).
- 746 Kavetski, D., and M.P. Clark (2010), Ancient numerical daemons of conceptual hydrolog-
747 ical modeling: 2. Impact of time stepping schemes on model analysis and prediction,
748 *Water Resources Research*, *46*, W10511, doi:10.1029/2009WR008896.
- 749 Laloy, E., and J.A. Vrugt (2012), High-dimensional posterior exploration of hydrologic
750 models using multiple-try DREAM_(ZS) and high-performance computing, *Water Re-*
751 *sources Research*, *48*, W01526, doi:10.1029/2011WR010608.
- 752 Lochbühler, T., J.A. Vrugt, M. Sadegh, and N. Linde (2014), Summary Statistics from
753 Training Images as Prior Information in Probabilistic Inversion, *Geophysical Research*
754 *Letters*, Under Review.
- 755 Marjoram, P., J. Molitor, V. Plagnol, and S. Tavaré (2003), Markov chain Monte Carlo
756 without likelihoods, *Proceedings of the National Academy of Sciences of the United*
757 *States of America*, *100* (26), 15324–15328.

- 758 Metropolis, N., A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller (1953),
759 Equation of state calculations by fast computing machines, *Journal of Chemical Physics*,
760 *21*, 1087–1092.
- 761 Nott, D.J., L. Marshall, and J. Brown (2012), Generalized likelihood uncertainty estima-
762 tion (GLUE) and approximate Bayesian computation: What’s the connection?, *Water*
763 *Resources Research*, *48*(12), doi:10.1029/2011WR011128
- 764 Price, K.V., Storn, R.M., and J.A. Lampinen (2005), *Differential evolution, A practical*
765 *approach to global optimization*, Springer, Berlin.
- 766 Pritchard, J.K., M.T. Seielstad, A. Perez-Lezaun, and M.T. Feldman (1999), Popula-
767 tion Growth of Human Y Chromosomes: A Study of Y Chromosome Microsatellites,
768 *Molecular Biology and Evolution*, *16*(12), 1791–1798.
- 769 Ratmann, O., C. Andrieu, C. Wiuf, and S. Richardson (2009), Model criticism based on
770 likelihood-free inference, with an application to protein network evolution, *Proceedings*
771 *of the National Academy of Sciences of the United States of America*, *106*, 1–6.
- 772 Sadegh, M., and J.A. Vrugt (2013), Approximate Bayesian Computation in hydrologic
773 modeling: equifinality of formal and informal approaches, *Hydrology and Earth System*
774 *Sciences - Discussions*, *10*, 4739–4797, doi:10.5194/hessd-10-4739-2013.
- 775 Schoups, G., and J.A. Vrugt (2010), A formal likelihood function for parameter and pre-
776 dictive inference of hydrologic models with correlated, heteroscedastic and non-Gaussian
777 errors, *Water Resources Research*, *46*, W10531, doi:10.1029/2009WR008933
- 778 Sisson, S.A., Y. Fan, and M.M. Tanaka (2007), Sequential Monte Carlo without likeli-
779 hoods, *Proceedings of the National Academy of Sciences of the United States of America*,
780 *104*(6), 1760–1765.

- 781 Smith, T., A. Sharma, L. Marshall, R Mehrotra, and S. Sisson (2010), Development of
782 a formal likelihood function for improved Bayesian inference of ephemeral catchments,
783 *Water Resources Research*, *46*(12), W12551.
- 784 Sunnåker M, A.G. Busetto, E. Numminen J. Corander M. Foll, and C. Dessimoz (2013),
785 Approximate Bayesian Computation, *Plos Computational Biology*, *9*(1), e1002803, 1–
786 10, doi:10.1371/journal.pcbi.1002803.
- 787 Storn, R., and K. Price (1997), Differential evolution - a simple and efficient heuristic for
788 global optimization over continuous spaces, *Journal of Global Optimization*, *11*, 341–
789 359.
- 790 Toni, T., D. Welch, N. Strelkowa, A. Ipsen, and M.P. Stumpf (2009), Approximate
791 Bayesian computation scheme for parameter inference and model selection in dynamical
792 systems, *Journal of the Royal Society Interface*, *6*, 187–202.
- 793 Turner, B.M, and T. van Zandt (2012), A tutorial on approximate Bayesian computation,
794 *Journal of Mathematical Psychology*, *56*, 69–85.
- 795 Turner, B.M, and P.B. Sederberg (2012), Approximate Bayesian computation
796 with differential evolution, *Journal of Mathematical Psychology*, *56*(5), 375–385,
797 doi:10.1016/j.jmp.2012.06.004.
- 798 Vrugt, J.A., C.J.F. ter Braak, M.P. Clark, J.M. Hyman, and B.A. Robinson (2008),
799 Treatment of input uncertainty in hydrologic modeling: Doing hydrology backward
800 with Markov chain Monte Carlo simulation, *Water Resources Research*, *44*, W00B09,
801 doi:10.1029/2007WR006720.
- 802 Vrugt, J.A., C.J.F. ter Braak, C.G.H. Diks, D. Higdon, B.A. Robinson, and J.M. Hyman
803 (2009), Accelerating Markov chain Monte Carlo simulation by differential evolution

804 with self-adaptive randomized subspace sampling, *International Journal of Nonlinear*
805 *Sciences and Numerical Simulation*, 10(3), 273–290.

806 Vrugt, J.A., and C.J.F. ter Braak, (2011), DREAM_(D): an adaptive Markov Chain Monte
807 Carlo simulation algorithm to solve discrete, noncontinuous, and combinatorial posterior
808 parameter estimation problems, *Hydrology and Earth System Sciences*, 15, 3701–3713,
809 doi:10.5194/hess-15-3701-2011.

810 Vrugt, J.A., and M. Sadegh (2013), Toward diagnostic model calibration and
811 evaluation: Approximate Bayesian computation, *Water Resources Research*, 49,
812 doi:10.1002/wrcr.20354.

813 Wegmann, D., C. Leuenberger, and L. Excoffier (2009), Efficient approximate Bayesian
814 computation coupled with Markov chain Monte Carlo without likelihood, *Genetics*,
815 182(4), 1207–1218.

Table 1: Case study I: 1-dimensional toy example. We list the (final) epsilon value, acceptance rate, AR (%) and number of function evaluations needed to sample the target distribution.

	ϵ	AR(%)	Func. Eval.
ABC-REJ	0.025	0.259	386,742
ABC-PMC	0.025	0.585	170,848
DREAM _(ABC)	0.025	1.708	50,000

Table 2: Case study II: 20-dimensional bivariate Gaussian distribution. We list the final epsilon value, acceptance rate, AR (%) and number of function evaluations needed to sample a stationary distribution.

	ϵ	AR(%)	Func. Eval.
ABC-PMC	0.600	0.019	5,143,989
DREAM _(ABC)	0.025	19.717	200,000

Table 3: Prior ranges of the SAC-SMA model parameters and their posterior mean values derived for the French Broad river basin data using the ABC-PMC and DREAM_(ABC) algorithms. We also summarize the results of a residual-based Gaussian likelihood function, DREAM_(RBGL).

Parameter	Range	Posterior Parameter Mean		
		ABC-PMC	DREAM _(ABC)	DREAM _(RBGL)
UZWWM	1-150	81.812	82.676	39.725
UZFWM	1-150	73.722	67.945	10.531
UZK	0.1-0.5	0.299	0.301	0.387
PCTIM	0-0.1	0.007	0.010	0.003
ADIMP	0-0.4	0.118	0.121	0.189
ZPERC	1-250	132.350	130.883	109.207
REXP	1-5	2.972	2.847	4.880
LZWWM	1-500	396.792	363.329	456.619
LZFSM	1-1000	282.049	246.326	146.584
LZFPM	1-1000	796.113	728.855	733.093
LZSK	0.01-0.25	0.145	0.142	0.124
LZPK	0.0001-0.025	0.006	0.006	0.009
PFREE	0-0.6	0.292	0.304	0.530
RRC	0-1	0.630	0.581	0.321

Table 4: Case study III: 14-dimensional SAC-SMA model calibration problem. We list the final epsilon value, acceptance rate, AR (%) and number of function evaluations needed to sample the posterior distribution. We also include the results of DREAM using a residual-based Gaussian likelihood function.

	ϵ	AR(%)	Func. Eval.
ABC-PMC	0.025	0.046	2,173,490
DREAM _(ABC)	0.025	3.135	200,000
DREAM _(RGL)	N/A	4.543	200,000

Table 5: Performance of the SAC-SMA model for the calibration and evaluation data period of the French Broad river basin. We summarize the coverage (%) and average width (mm/d) of the 95% simulation intervals (due to parameter uncertainty), and the RMSE (mm/d), bias (%) and correlation coefficient, R (-) of the posterior mean SAC-SMA simulation.

	Coverage (%)		Width (mm/d)		RMSE (mm/d)		Bias (%)		R (-)	
	Calib.	Eval.	Calib.	Eval.	Calib.	Eval.	Calib.	Eval.	Calib.	Eval.
ABC-PMC	70.991	65.328	1.301	1.185	0.932	0.812	3.281	-3.721	0.863	0.832
DREAM _(ABC)	71.100	66.515	1.408	1.272	0.831	0.722	3.731	-3.245	0.893	0.866
DREAM _(RBGL)	18.829	12.500	0.155	0.141	0.539	0.545	2.767	-0.042	0.956	0.924

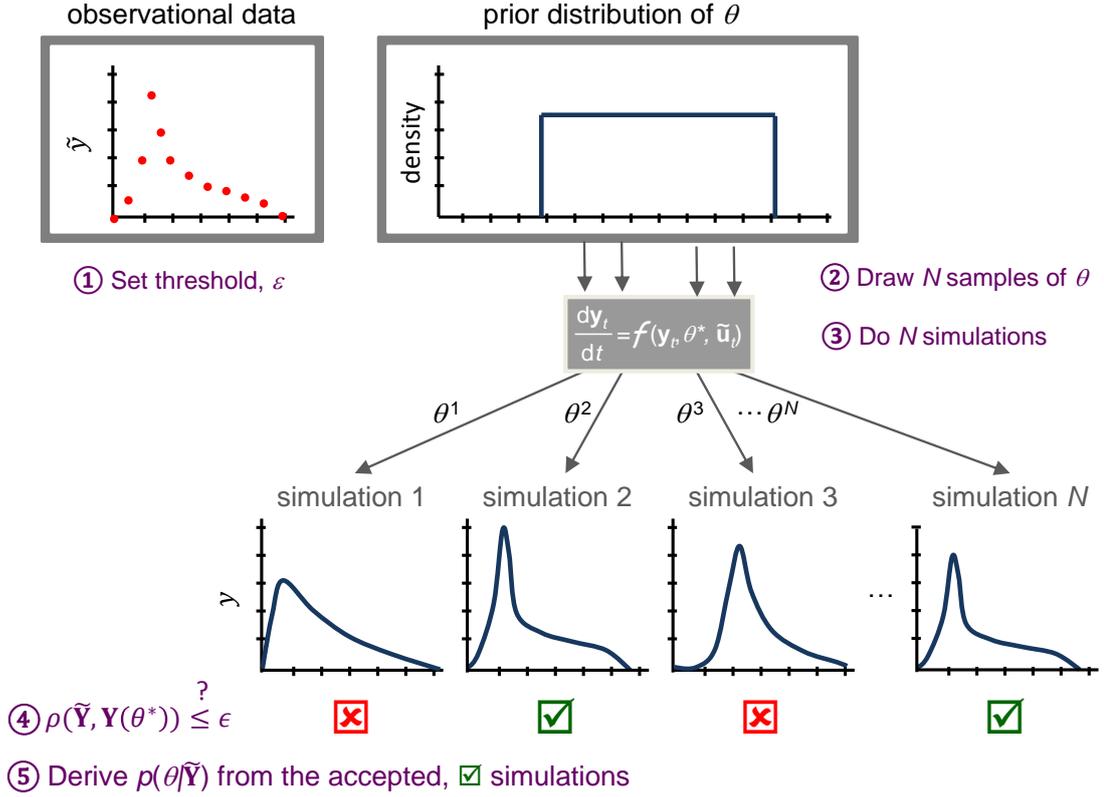


Figure 1: Conceptual overview of approximate Bayesian computation (ABC) for a hypothetical one-dimensional parameter estimation problem. First, N samples are drawn from a user-defined prior distribution, $\theta^* \sim p(\theta)$. Then, this ensemble is evaluated with the model and creates N model simulations. If the distance between the observed and simulated data, $\rho(\tilde{Y}, Y(\theta^*))$ is smaller than or equal to some nominal value, ϵ then θ^* is retained, otherwise the simulation is discarded. The accepted samples are then used to approximate the posterior parameter distribution, $p(\theta|\tilde{Y})$. Note that for sufficiently complex models and large data sets the probability of happening upon a simulation run that yields precisely the same simulated values as the observations will be very small, often unacceptably so. Therefore, $\rho(\tilde{Y}, Y(\theta^*))$ is typically defined as a distance between summary statistics of the simulated, $S(Y(\theta^*))$ and observed, $S(\tilde{Y})$ data, respectively. Modified after *Sunnåker et al.* [2013].

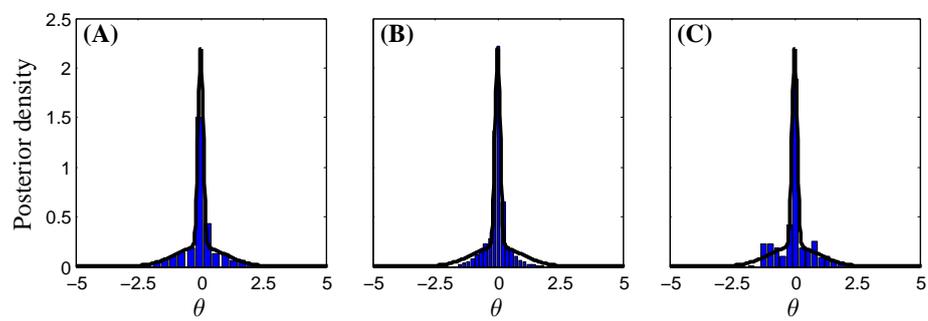


Figure 2: One-dimensional mixture distribution (solid black line) and histogram of the posterior samples derived from (A) Rejection sampling (ABC-REJ), (B) Population Monte Carlo sampling (ABC-PMC), and (C) MCMC simulation with DREAM_(ABC).

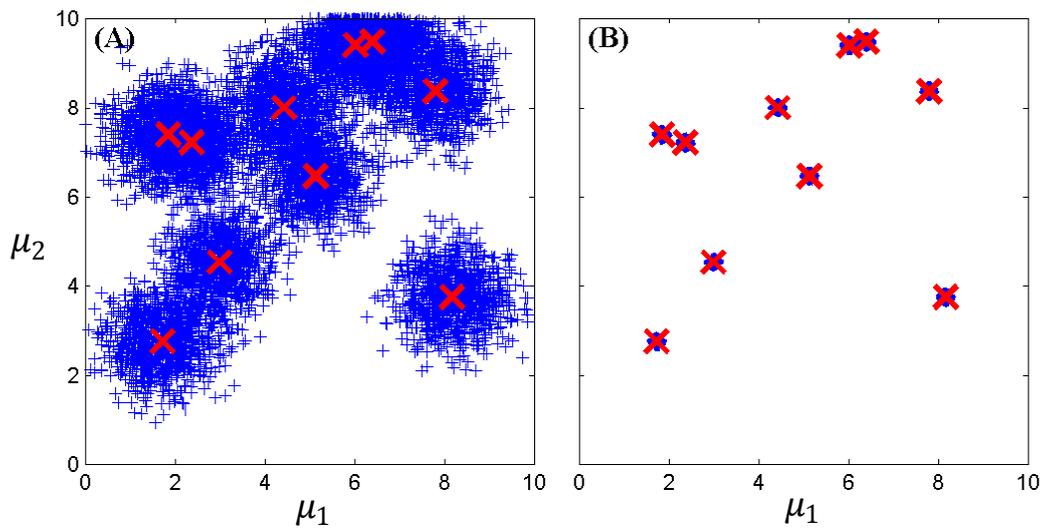


Figure 3: Two-dimensional scatter plots of the posterior samples, "+" generated with (A) Population Monte Carlo sampling, and (B) MCMC simulation with DREAM_(ABC). The true values of μ_1 and μ_2 are indicated with a cross, "x". The ABC-PMC samples are very dispersed and significantly over estimate the actual width of the target distribution.

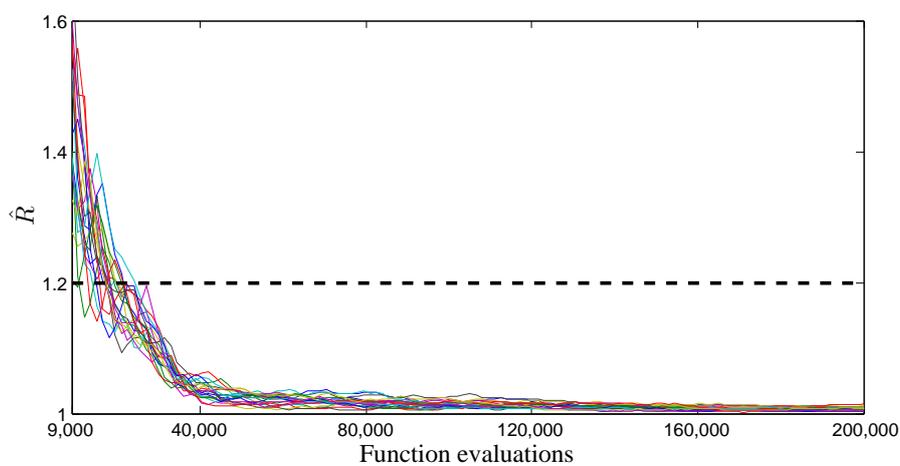


Figure 4: Trace plots of the \hat{R} -statistic of the sampled Markov chains with DREAM_(ABC) for the 20-dimensional bivariate normal distribution. Each parameter is coded with a different color. The dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

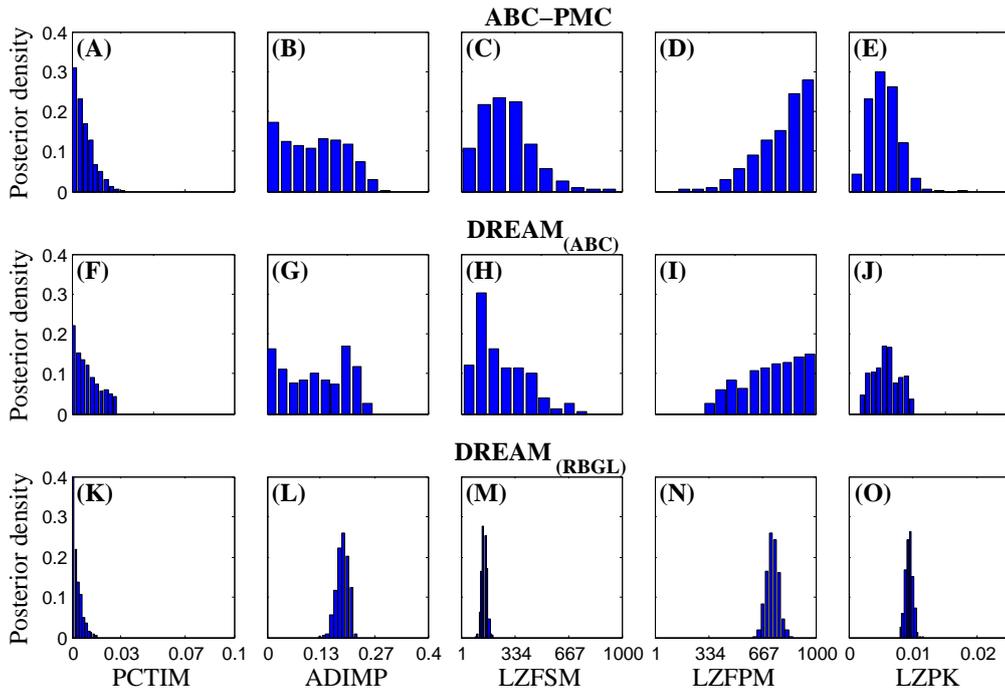


Figure 5: Marginal distributions of PCTIM, ADIMP, LZFSM, LZFPM and LZPK derived from the posterior samples created with ABC-PMC (upper panel) and DREAM_(ABC) (middle panel). The histograms of the five SAC-SMA parameters occupy almost their entire prior distribution, which suggests that they are not particularly well identifiable by calibration against the observed annual baseflow index, annual runoff coefficient, and flow duration curve, respectively. The results of both ABC sampling methods are in good (visual) agreement, which inspires confidence in the ability of DREAM_(ABC) to correctly sample the underlying target distribution. The bottom panel displays histograms of the SAC-SMA parameters derived using a classical residual-based (Gaussian) likelihood function. The SAC-SMA parameters are much better resolved, but these results cannot be justified given (amongst others) a lack of treatment of rainfall data errors.

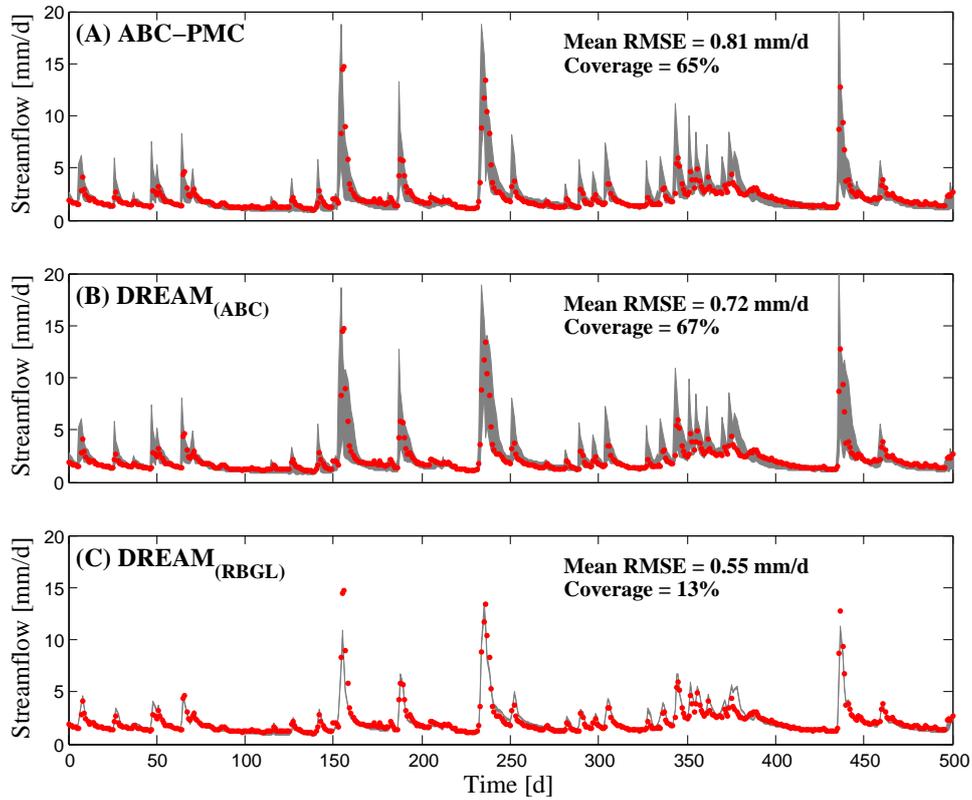


Figure 6: 95% posterior simulation uncertainty ranges (gray region) of the SAC-SMA model for a selected portion of the evaluation data set of the French Broad watershed. The top two panels displays the results of diagnostic model evaluation using (A) ABC-PMC, and (B) DREAM_(ABC), whereas the bottom panel depicts the results of DREAM with a classical residual-based Gaussian likelihood function. The observed discharge values are indicated with the red dots. The SAC-SMA simulation intervals derived from ABC-PMC and DREAM_(ABC) are very similar and encapsulate a large part of the discharge observations. The DREAM_(RBGL) uncertainty ranges, on the other hand, exhibit a much lower coverage, but closer track the observed discharge data.

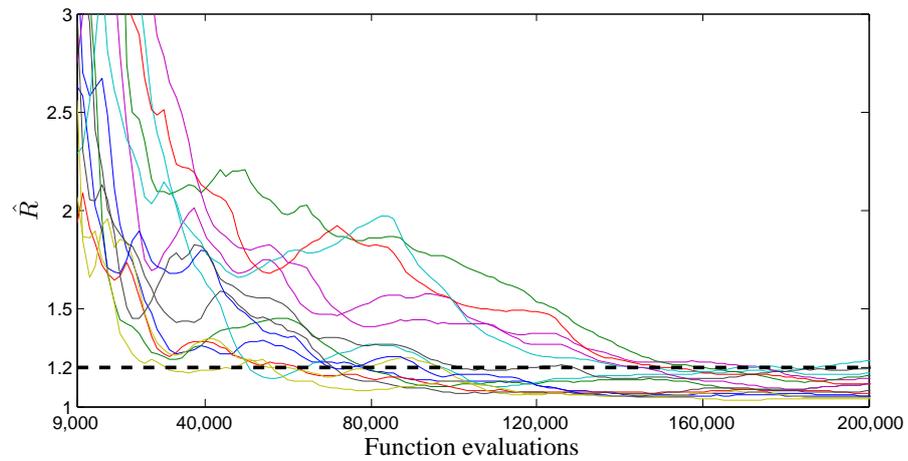


Figure 7: Evolution of the \hat{R} -statistic for the parameters of the SAC-SMA model using DREAM_(ABC) and discharge data from the French Broad watershed. Each of the parameters is coded with a different color. The dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

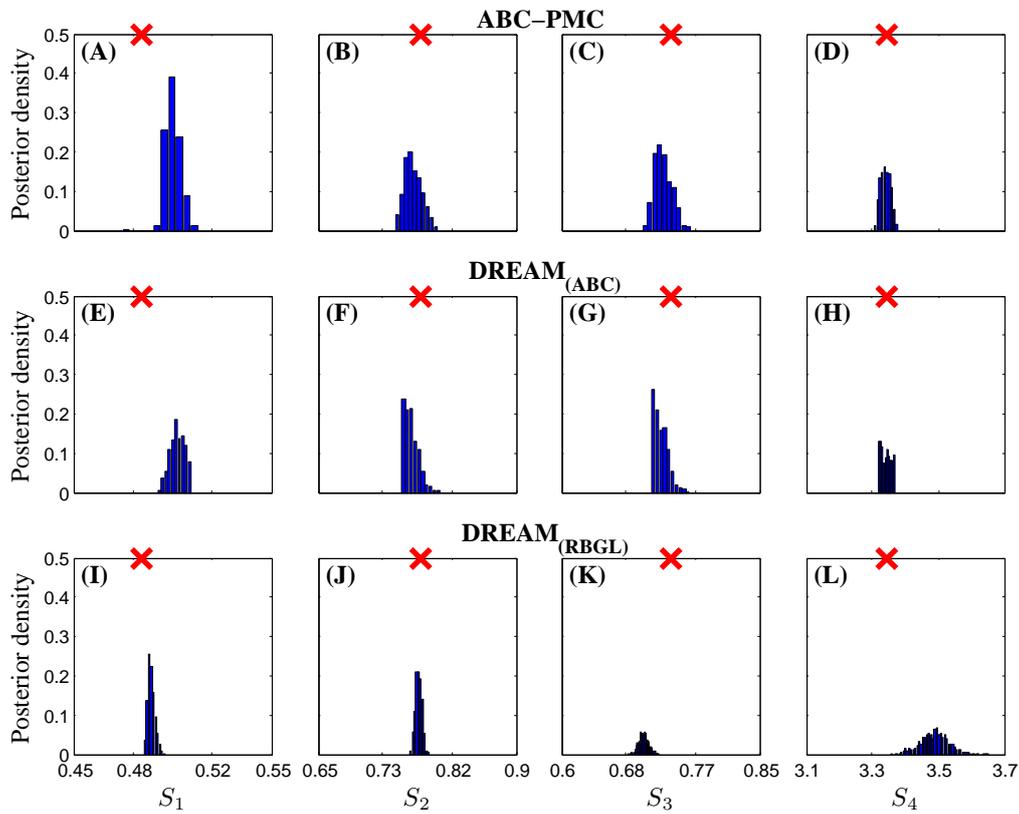


Figure 8: Histograms of the SAC-SMA derived summary statistics, S_1 (runoff), S_2 (baseflow), S_3 and S_4 (flow duration curve) of the posterior samples from ABC-PMC (top panel) and DREAM_(ABC) (middle panel). The bottom panel displays the results of DREAM with a residual-based Gaussian likelihood function. The observed values of the summary metrics are separately indicated in each plot using the "x" symbol. While $S_2 \rightarrow S_4$ center around their observed value (x) for the ABC analysis, the marginal posterior distribution of S_1 is skewed to the right and does not encapsulate its measured value. This demonstrates that model is unable to simultaneously satisfy all the four different summary metrics used herein. This

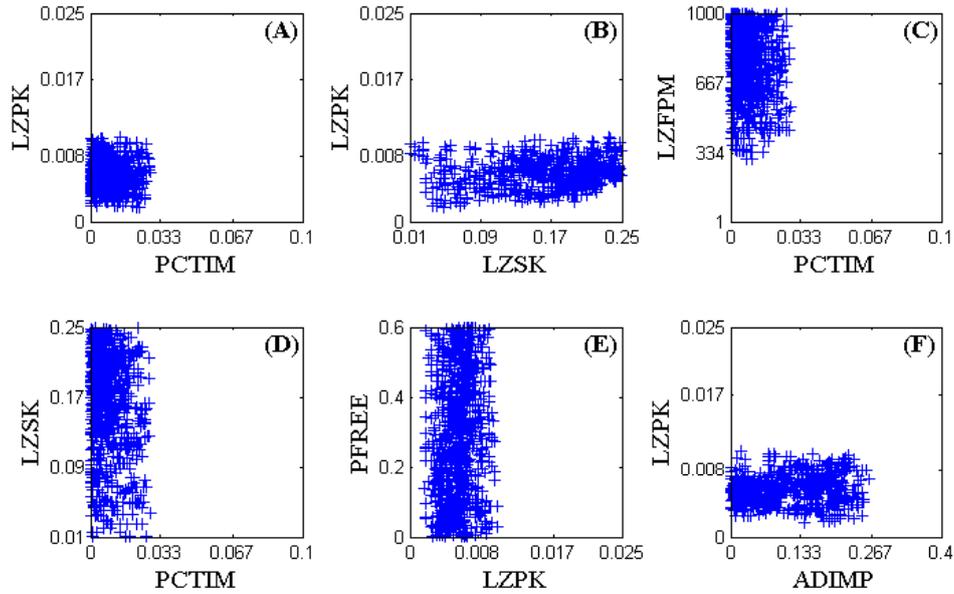


Figure 9: Two-dimensional scatter plots of the posterior samples, ”+” generated with DREAM_(ABC) for six different pairs of parameters of the SAC-SMA model. We restrict our attention to the (A) PCTIM - LZPK, (B) LZSK - LZPK, (C) PCTIM - LZFPM, (D) PCTIM - LZSK, (E) LZPK - PFREE, and (F) ADIMP - LZPK space, respectively. The bivariate posterior samples occupy a well-defined hypercube interior to the prior distribution.