

24 research to define an alternative and more practical convergence criteria for MCMC applications to
25 computationally intensive hydrologic models may be warranted.

26 *Keywords:* Hydrologic modelling, Multi-criteria calibration, Uncertainty analysis, Bayesian inference,
27 GLUE

28 **1 Introduction**

29 Hydrologic modelling has benefited from significant developments over the past two
30 decades, and this has led to increasing complexity in hydrologic models and an advance from
31 lumped conceptual models toward semi-distributed and distributed physics-based models. These
32 models include many parameters which need to be estimated through an adjustment procedure
33 using historical observation data. The automatic calibration conducted without sufficient
34 hydrological expertise might yield improper parameter values which can result in unreasonable
35 regimes of model responses that are not controlled by measurements (Refsgaard, 1997; Wagener
36 et al., 2001). Moreover, even ‘well calibrated’ parameter values can yield poor performance with
37 respect to an independent validation data set.

38 Problems with parameter adjustment in hydrologic models can be attributed to different
39 factors. Conceptually, aggregation of all residuals into a single objective function during
40 calibration does not provide sufficient detail about model inadequacy (Gupta et al., 1998). For
41 example, single-objective metrics do not distinguish between high-flow and low-flow model
42 behaviour. This realization has motivated multi-criteria calibration approaches in which multiple
43 sets of observations and/or multiple evaluation criteria are employed (Gupta et al., 1998; Legates
44 and McCabe, 1999; Madsen, 2000; Yapo et al., 1998). Multi-criteria calibration uses more than
45 one index to describe the characteristics of the error vector (e.g., separate Nash-Sutcliffe values

46 for high-flow and low-flow data), resulting in an objective-function tradeoff curve and
47 corresponding set of “Pareto” optimal parameter values.

48 Another strategy for increasing the usefulness of predictive hydrologic models is to
49 rigorously account for different sources of uncertainty (e.g., uncertainties associated with
50 estimated parameter values as well as uncertainties in meteorological inputs and other non-
51 calibrated forcing functions). In fact, it is very important to include an assessment of uncertainty
52 in the calibration process. Razavi et al. (2010) named such approaches ‘uncertainty-based
53 calibration’ which refers to the coupling of an environmental model with an uncertainty engine
54 such that the uncertainty engine repeatedly samples model parameter configurations to develop a
55 calibrated probability distribution for the parameters. Other research has emphasized
56 comprehensive model assessment (or model evaluation) procedures whereby parameter
57 estimation is done probabilistically to derive the probability density function (PDF) of the model
58 outcome(s) of interest, through traditional ‘frequentist’ approaches (e.g. Bates and Watts, 1988;
59 Reichert, 1997; Seber and Wild, 1989) and Bayesian inference approaches.

60 From a Bayesian perspective, uncertainty-based calibration seeks to elucidate posterior PDFs
61 for various parameters and model outcomes given some prior information and available data.
62 These posterior PDFs then form the basis of a complementary predictive uncertainty analysis
63 (Bates and Campbell, 2001; Box and Tiao, 1973; Gelman et al., 2004; Kavetski et al., 2002;
64 Kuczera, 1983; Kuczera and Parent, 1998; Thiemann et al., 2001). The Bayesian approach to
65 model specification and uncertainty analysis is particularly appealing as it allows for formal
66 specification and propagation of an error model (Marshall et al., 2007). Furthermore, in the
67 Bayesian approach, any a priori knowledge about model parameters can be used in terms of prior
68 distributions, which are then updated for any particular catchment using the data available. For

69 complex hydrologic models, Bayesian inference is aided by the use of numerical procedures that
70 implement Markov Chain Monte Carlo (MCMC) sampling. In this regard, a number of MCMC
71 samplers have been proposed, including BaRE (Bayesian Recursive Estimation) (Thiemann et
72 al., 2001); SCEM-UA (Shuffled Complex Evolution Metropolis – University of Arizona) (Vrugt
73 et al., 2003b), BATEA (Bayesian Total Error Analysis) (Kavetski et al., 2002; Kavetski et al.,
74 2006), and DREAM (Differential Evolution Adaptive Metropolis) (Vrugt et al., 2009).

75 At the heart of Bayesian inference is the use of formal likelihood functions to analyse
76 parameter uncertainty. A given likelihood function must make explicit assumptions about the
77 form of the model residuals (i.e., deviations between simulations and observations) (Stedinger et
78 al., 2008). Thus, a major criticism of the Bayesian approach is that in hydrologic modelling the
79 appropriate statistical form for a given set of model residuals is not always clear, and this makes
80 it difficult to establish an appropriate likelihood function (e.g., Beven et al., 2008). To address
81 this issue, some researchers have emphasized the development of more appropriate likelihood
82 functions by using hierarchical Bayesian structures that disaggregate different sources of
83 uncertainties (e.g., Huard, 2008; Kuczera et al., 2006; Moradkhani et al., 2005; Renard et al.,
84 2010; 2011; Wei et al., 2010). However, development and application of such formulations to
85 complex non-linear hydrological models is non-trivial and may be computationally intractable in
86 some case studies using existing state-of-the-art MCMC samplers. The issue of defining an
87 appropriate Bayesian likelihood formulation becomes even more challenging when one considers
88 a multi-response or multi-criteria approach – an approach that some have argued is the most
89 appropriate for hydrological modelling (e.g., Hamilton, 2007; Montanari, 2007).

90 Recently, the concept of epistemic and aleatory uncertainties in hydrological modelling has
91 been discussed among researchers (Beven et al., 2012; Beven et al., 2011; Clark et al., 2012;

92 Montanari, 2011). Uncertainties are categorized as aleatory (also called natural uncertainty) if
93 they are presumed to be the intrinsic randomness of a stochastic process which can be
94 represented in terms of the probabilities of different outcomes. On the other hand, many of the
95 errors that enter into the modelling process stem from a lack of knowledge about processes and
96 boundary conditions. These errors are called epistemic or limited-knowledge uncertainty. In
97 statistical models (including Bayesian inference structures), uncertainties are accounted for by
98 providing a representation of all of the important sources of uncertainty as aleatory (Beven et al.,
99 2011). As a consequence, the results of Bayesian methods might not be robust when many of the
100 errors that affect modelling uncertainty in hydrology are epistemic (Beven et al., 2011; Beven et
101 al., 2008). However, statistical methods are believed to be able to fit epistemic uncertainties
102 provided that the inherent regularities are well represented by the statistical model itself
103 (Montanari, 2011). Similar to almost all studies in the literature on uncertainty analysis of
104 rainfall-runoff models, the Bayesian method of our paper also considers all uncertainties to be
105 aleatory.

106 Despite the robust theoretical underpinnings of a formal Bayesian approach to parameter
107 inference, a variety of alternative and informal approaches have been proposed for uncertainty-
108 based multi-criteria calibration of complex hydrological models. Examples include a Pareto-
109 based calibration approach (Gupta et al., 1998) and informal MCMC sampling (Blasone et al.,
110 2008a; Vrugt et al., 2003a). Importance sampling techniques have also been used for informal
111 uncertainty-based calibration, with GLUE (Generalized Likelihood Uncertainty Estimation)
112 (Beven and Binley, 1992) being the most commonly used approach. GLUE is based on the
113 concept of ‘equifinality’ and classifies Monte Carlo samples as having produced model output
114 that is either ‘behavioural’ (i.e., plausible, given the data and one’s knowledge of the system) or

115 'non-behavioural'. The behavioural solutions are then used to derive the probability distribution
116 function for parameters and model outputs. The GLUE methodology can be easily extended to
117 multi-criteria calibration problems (e.g. Blazkova and Beven, 2009). A drawback of informal
118 methods is that such approaches do not require formal specification of an error model and might
119 not be reliable for uncertainty analysis (Kavetski et al., 2002).

120 Along with development of a variety of uncertainty-based calibration routines, some
121 researchers have focused their efforts towards comparison between formal and informal
122 methods. Overall, these efforts generally indicate relatively close agreement among alternative
123 methods, in terms of predictive capability (Beven et al., 2008; Jeremiah et al., 2011; Jin et al.,
124 2010; Li et al., 2010; Qian et al., 2003; Vrugt et al., 2008; Yang et al., 2008). Note that some
125 studies have only considered informal methods in their comparisons (e.g., Blasone et al., 2008b).

126 From both a comparative and theoretical perspective, previous literature demonstrates that
127 MCMC sampling and Bayesian inference can be considered a preferred approach to deal with
128 uncertainty-based calibration, as long as the computational budget allows full convergence of the
129 MCMC sampler. Achieving convergence is not problematic if one is dealing with rainfall-runoff
130 models with manageable simulation runtimes. However, when computational budget limitations
131 exist, MCMC sampling may not be an appropriate choice. Furthermore, the observed similarity
132 between the predictive capabilities of formal and informal approaches suggests that one might be
133 able to gain insight into predictive uncertainty by means of informal approaches without getting
134 involved in likelihood definition and corresponding assumptions. Most of previous papers
135 comparing formal and informal approaches have only considered single-criterion calibration
136 scenarios. Balin-Talamba (2004) and Balin-Talamba et al. (2010) considered multi-criteria
137 calibration of hydrologic models applying GLUE and MCMC sampling. These studies evaluated

138 the impact of multi-response calibration on predictive uncertainty using GLUE and MCMC, in
139 comparison with single-criterion calibration. However, the GLUE and MCMC techniques are
140 only visually compared in Balin-Talamba (2004) and no comparative measures are reported. To
141 the best of our knowledge, comparison among formal and informal techniques from a multi-
142 criteria perspective using quantitative comparative measures has yet to be reported on in the
143 literature.

144 The main objective of this research is to evaluate the applicability of different uncertainty
145 analysis approaches to multi-criteria calibration and uncertainty analysis of hydrologic models
146 considering identical computational budget. The methodologies addressed in this paper are
147 statistically-based Bayesian inference using MCMC sampling (Bates and Campbell, 2001;
148 Kuczera, 1983; Schaefli et al., 2007; Vrugt et al., 2009), and sampling-based uncertainty
149 estimation using GLUE (Beven and Binley, 1992; Blazkova and Beven, 2009). Bayesian
150 inference was implemented using the DREAM MCMC sampler (Vrugt et al., 2009) through a
151 robust multi-criteria formulation. Also, we consider an alternative Bayesian method based on the
152 results of MCMC sampling up to a limited computational budget (i.e., using the MCMC before
153 convergence). Such a method cannot be viewed informal, as it uses formal likelihood function;
154 however, it would not be formal either, as convergence has not occurred, meaning that the
155 solutions in the chain could not be considered as samples from posterior distributions.

156 **2 Methodology**

157 A typical multi-criteria model calibration process can involve multiple likelihood functions
158 used for different sets of measurements, e.g., discharge, sediment, snow, etc. However, even in
159 the case of a model with only one output flux to be simulated, the model evaluation may still be

160 considered to be inherently multi-criteria (Gupta et al., 1998). The multi-criteria numerical
161 experiments in this study only deal with one response (discharge), splitting it into high- and low-
162 flows. This strategy is expected to be adequate for an initial exploration of multiple uncertainty-
163 based calibration techniques within a multi-criteria formulation.

164 The comparison framework of this study uses the posterior distribution of model parameters
165 derived from MCMC sampling, as well as the behavioural or optimal parameter sets obtained
166 from other methods. In order to be consistent in wording, the term “posterior” is applied to all of
167 the considered techniques even though the results of non-converged MCMC sampling and
168 GLUE are not a formal statistical posterior distribution. Results are then compared with respect
169 to computational burden, complexity, and predictive capacity. Numerical experiments are aimed
170 at exploring advantages and disadvantages of the uncertainty analysis techniques addressed in
171 this study in multi-criteria calibration of rainfall-runoff models. The reliability of these methods
172 is evaluated using two rainfall-runoff models, a 5-parameter lumped model, HYMOD
173 (Hydrology model) (Boyle, 2000), and an 11-parameter semi-distributed model, WetSpa (Water
174 and Energy Transfer between Soil, Plants and Atmosphere) (Liu et al., 2003; Wang et al., 1996).

175 The GLUE approach of this paper employs informal likelihood functions and results are
176 compared with those obtained from formal Bayesian inference as well as non-converged MCMC
177 sampling. The use of GLUE without a formal likelihood function has been the subject of much
178 debate (e.g., Beven et al., 2008; Mantovan and Todini, 2006; Montanari, 2005; Thiemann et al.,
179 2001). Nevertheless, we used GLUE with an informal generalized likelihood function in this
180 study because the objective of the study was to assess the performance of informal methods.
181 Much of the reason informal methods like GLUE are so well utilized in practice is because they
182 can use informal likelihood functions based on long utilized deterministic calibration objective

183 functions like sum of squared errors or the Nash Sutcliffe coefficient. It is also worth noting that
184 GLUE could also be applied using formal likelihood functions (Freni and Mannina, 2009;
185 Romanowicz et al., 1994), but this is not addressed in the present paper.

186 The comparison approach (informal to formal methods) of this study is exactly consistent
187 with previous comparative studies of uncertainty-based calibration in hydrological modelling
188 (e.g. Vrugt et al., 2008; Yang et al., 2008). Beven (2009) noticed that in Vrugt et al. (2008) the
189 formal Bayes estimates are based on an autoregressive error model, while such information is not
190 supplied to the GLUE simulations. Despite the difference between the formulations of the
191 Bayesian approach and GLUE in Vrugt et al. (2008), it is shown in that paper that formal and
192 informal uncertainty analysis methods have some common ground with respect to the total
193 predictive uncertainty in single-criterion calibration cases. In this paper, multiple quantitative
194 comparative measures are applied and we evaluate the similarity in behavior of MCMC and
195 GLUE in the multi-criteria context. As such, we consider the same implementations of MCMC
196 sampling and GLUE as used in Vrugt et al. (2008).

197 **2.1 Formal multi-criteria Bayesian inference**

198 Bayesian statistics have been shown to be a robust methodology for formal multi-criteria
199 calibration and uncertainty analysis of hydrologic models, as long as all underlying assumptions
200 are satisfied. Both analytical and numerical Bayesian approaches have been used to deal with
201 multi-criteria calibration (Balin-Talamba et al., 2010; Hong et al., 2005; Kuczera, 1983; Kuczera
202 and Mroczkowski, 1998; Mroczkowski et al., 1997; Schaeffli et al., 2007). The notion of multi-
203 criteria in Bayesian inference structures is mostly concerning cases in which multiple responses
204 of observations are employed (e.g., measured streamflows and measured soil water content), and

205 thus, it is also called multi-response calibration in the literature. There are also reports of multi-
206 criteria Bayesian formulations using a single response. For instance, Schaepli et al. (2007)
207 considered multiple likelihood functions which were associated with high- and low-streamflows.
208 The research presented here used a previously published multi-criteria formulation (Balin-
209 Talamba et al., 2010; Schaepli et al., 2007).

210 Moreover, the initial experiments of Bayesian inference in these case studies showed that
211 errors were correlated. As a result, we had to consider development of a formal likelihood
212 function which accounts for auto-correlation. As such, auto-regressive (AR) parameters were
213 introduced to the high- and low-flow time series to address auto-correlation among residuals
214 (e.g., Bates and Campbell, 2001; Kuczera, 1983). The resulting Bayesian inference formulation
215 introduces a first-order AR scheme to represent the residuals (Balin-Talamba et al., 2010;
216 Schaepli et al., 2007), details of which are provided in the Appendix of this paper. Note that the
217 AR scheme was applied separately to the low- and high-flow regimes and this resulted in the
218 addition of two AR parameters (ρ_L for low-flows and ρ_H for high-flows) to the set of calibrated
219 parameters.

220 For this paper, the DREAM MCMC sampler was used for formal Bayesian inference (Vrugt
221 et al., 2009). DREAM maintains ergodicity while showing excellent efficiency even if the target
222 posterior distributions are complex, highly nonlinear, and/or multimodal. DREAM runs multiple
223 Markov chains simultaneously to facilitate efficient global exploration of the parameter space.
224 Like other adaptive samplers, DREAM speeds convergence by dynamically adjusting the scale
225 and orientation of the proposal distribution.

226

227 **2.2 Sampling-based Uncertainty Estimation using Non-converged MCMC**

228 Even though applications of MCMC sampling with pseudo-likelihood functions have been
229 previously reported in the literature (Blasone et al., 2008b; Vrugt et al., 2003a), there has been no
230 report on evaluation of the results from non-converged MCMC samplers with formal likelihood
231 functions. In this paper, non-converged DREAM results are used to approximate the converged
232 MCMC sampling strategy. The number of solutions taken from a given DREAM chain was
233 defined to be consistent with the informal methods considered in this paper (explained below).
234 For example, if the informal methods use a budget of 10000 simulations, then we only consider
235 10000 solutions from the initial part of the long DREAM chain. Afterwards, the last 1000
236 solutions of this set would be treated as posterior solutions to derive prediction intervals. Clearly,
237 such an approach is neither formal (as convergence has not occurred) nor informal (as it uses
238 formal likelihood function). That is the reason why we separated this approach from formal
239 Bayesian and informal GLUE approaches.

240 **2.3 Sampling-based Uncertainty Estimation using GLUE**

241 The GLUE technique (Beven and Binley, 1992) is the most commonly applied method in the
242 family of informal sampling-based methods. In GLUE, parameter uncertainty accounts for all
243 sources of uncertainty, because “the likelihood measure value is associated with a parameter set
244 and reflects all these sources of error and any effects of the covariation of parameter values on
245 model performance implicitly” (Beven and Freer, 2001). The GLUE analysis conducted here
246 consisted of the following four steps:

247 1. Defining the generalized informal likelihood measure $l(\theta)$. Generally, the measure $l(\theta)$ is
248 a pseudo-likelihood function which demonstrates the model performance for a particular

249 parameter set θ . In this study, we used the generalized likelihood function provided in previous
 250 multi-criteria GLUE studies (Balin-Talamba, 2004; Lamb et al., 1998) as follows:

$$251 \quad l(\theta) = \prod_{i=1}^M \exp\left(-W_i \frac{\sigma_{\varepsilon,i}^2}{\sigma_{o,i}^2}\right) \quad (1)$$

252 where W_i represents the weighting factor for criterion i (explained later), M is the number
 253 of criteria, $\sigma_{\varepsilon,i}^2$ and $\sigma_{o,i}^2$ are the variance of simulation errors and the variance of observed data,
 254 respectively, over the time window in which criterion i is calculated. The likelihood function
 255 $l(\theta)$ equals 1 if the observed and simulated data are the same for all criteria, and reduces
 256 towards zero as the similarity decreases. Note that, in the multi-criteria calibration problem of
 257 this paper, we calculate this likelihood function based on the information in high- and low-flow
 258 time periods ($M=2$).

259 2. After defining $l(\theta)$, a large number of parameter sets are randomly sampled from the prior
 260 distribution and each parameter set is assessed as either “behavioural” or “non-behavioural”
 261 through a comparison of the likelihood measure with a selected threshold value which is
 262 explained in details later in this section of the paper.

263 3. Each behavioural parameter set is given a likelihood weight according to
 264 $\varpi_i = l(\theta_i) / \sum_{k=1}^N l(\theta_k)$, where N is the number of behavioural parameter sets.

265 4. Finally, prediction uncertainty of streamflow is described by quantiles of the cumulative
 266 distribution realized from the weighted behavioural parameter sets, i.e., at each time step, the
 267 model outcome associated to behavioural solutions are identified and prediction intervals (for
 268 example 95% intervals) are constructed based on quantiles (such as 2.5 and 97.5 percentiles).

269 The behavioural threshold for the GLUE pseudo-likelihood function defines the boundary
 270 between behavioural and non-behavioural solutions. In this study, based on the strategy in Balin-
 271 Talamba (2004) and Lamb et al. (1998), we followed the same strategy (also described below) to
 272 filter out behavioural samples. Once samples are taken from prior distributions, the generalized
 273 likelihood function Eq. (1) is calculated considering high- and low-flow time periods whereby
 274 the weights are equal for both periods, i.e., $W_L = 0.5$ and $W_H = 0.5$ (note that L and H stand for
 275 low- and high-flows, respectively):

$$276 \quad l(\theta) = \exp\left(-W_L \frac{\sigma_{\varepsilon,L}^2}{\sigma_{o,L}^2}\right) \cdot \exp\left(-W_H \frac{\sigma_{\varepsilon,H}^2}{\sigma_{o,H}^2}\right) = \exp\left[-\left(W_L \frac{\sigma_{\varepsilon,L}^2}{\sigma_{o,L}^2} + W_H \frac{\sigma_{\varepsilon,H}^2}{\sigma_{o,H}^2}\right)\right] \quad (2)$$

277 Parameter sets are now sorted based on the combined criterion, and the top N samples are
 278 considered behavioural solutions. Identifying N is in fact a subjective decision in GLUE, and
 279 would probably affect the uncertainty bounds computed using the GLUE method. Among the
 280 traditional choices reported in literature is N being equal to the number of top 10% of solutions
 281 (Binley and Beven, 1991; Lamb et al., 1998) sampled from the prior distributions. However,
 282 Lamb et al. (1998) showed that relaxation of the rejection threshold to define a larger proportion
 283 of the total number of samples as behavioural would cause only slight modifications of
 284 uncertainty bounds. The reason for this insensitivity to the rejection threshold is that even after
 285 selecting a larger number of behavioural samples, the majority of samples would achieve only
 286 small likelihood values. Therefore, the predictions associated with these poor samples would fall
 287 within the tails of the cumulative distributions of model outcome. Given the rescaling stage in
 288 GLUE, these predictions would have little effect on the location of uncertainty bounds (Lamb et
 289 al., 1998). In this paper, we also considered the top 10% strategy to define behavioural samples.

290 **2.4 Comparison measures**

291 The main goal of calibration and uncertainty analysis is to assess models' predictive
292 capability. Therefore, in order to evaluate uncertainty-based calibration techniques, it seems
293 necessary that we focus more on the validation time period rather than the calibration period.
294 Nonetheless, a portion of our analysis examined differences between calibration and validation
295 results. The comparative measures are calculated based on the results obtained using the
296 posterior parameter sets. It should be noted that the parameter uncertainty is derived based on the
297 envelope of model outputs using the posterior parameter sets. Moreover, in order to derive the
298 predictive uncertainty, the entire set of posterior parameters is first used in simulation model to
299 derive the parameter uncertainty. Afterwards, error parameters are sampled to generate a
300 correlated residual time series which is then added to model outputs.

301 To evaluate the quality of resulting model outcomes, efficiency measures such as NS is used
302 to assess model performance. In the multi-criteria context of this paper, we illustrate the
303 scatterplot of posterior parameter sets in bi-criteria space (i.e., NS for high- and low-flows).

304 In addition, the generated model outcomes using the posterior solutions derived from
305 different techniques are used to derive the predictive uncertainty which can be assessed using a
306 variety of measures. Among the simplest measures for comparing alternative realizations of
307 predictive uncertainty are the reliability and sharpness measures (Yadav et al., 2007). For a given
308 prediction interval, the reliability measure is the percentage of discharge observations that are
309 captured by the prediction interval. Reliability values are calculated by counting the number of
310 times the observed streamflow falls within the prediction band, divided by the length of the time
311 series. Sharpness is a measure of the prediction intervals' width relative to the hydrograph
312 prediction bounds obtained from sampling prior feasible parameter ranges. If the posterior

313 prediction bounds for the hydrograph form a single line, sharpness would be 100%. Whereas
314 when the posterior prediction bounds are the same as those obtained using priori feasible
315 parameter ranges, sharpness would be 0% (clearly undesirable). Ideally, and for a given
316 prediction interval, the reliability should be equal to the desired interval percentage (i.e., 90% of
317 observations should be captured by a 90% prediction interval) and larger values of the
318 corresponding sharpness measure are better than smaller values.

319 The Bayesian posterior predictive p-value is another measure of the predictive capacity of
320 uncertainty-based calibration techniques (Gelman et al., 2004, pp. 162–163). The Bayesian p-
321 value is the probability that the model prediction at a particular time step could be more extreme
322 than the observed data at that same time step. Such values may be estimated by the proportion of
323 simulations for which the simulated value equals or exceeds the observed value. Probability
324 distributions of p-values can be constructed from the complete series of p-value calculations. If
325 the model output and measured data are consistent, the corresponding p-value distribution should
326 be uniformly distributed over the interval [0,1]. This can be checked graphically using QQ-plots
327 (Laio and Tamea, 2007; Thyer et al., 2009) and deviations from the bisector (the 1:1 line) denote
328 interpretable deficiencies (see Figure 1).

329 **[Figure 1 goes here.]**

330 Our approach to compute comparative performance metrics with GLUE such as reliability,
331 sharpness and Bayesian p-values is consistent with studies computing one or more of these
332 metrics for GLUE results based on a pseudo-likelihood function such as Vrugt et al. (2008),
333 Yang et al. (2008) and Jin et al. (2010)

334 2.5 Case Studies

335 Bayesian inference is expected to result in robust expression of predictive uncertainty, as
336 long as all assumptions are satisfied and the posterior PDFs are taken from a converged MCMC
337 sampler. Two case-studies involving real data from two catchments are used in this paper, for
338 which the DREAM sampler is run to convergence to extract formal posterior distributions. The
339 non-converged MCMC sampling and GLUE methods are also applied to the same problems. One
340 case-study applies the HYMOD hydrologic model to the Leaf River catchment, and one applies
341 the WetSpa hydrologic model to the Hornad River catchment, where details about these
342 catchments are provided below.

343 The first study area addressed in this paper is the 1994 km² Leaf River watershed located
344 north of Collins, Mississippi. This catchment has been studied intensively in the past (e.g.,
345 Boyle, 2000; Sorooshian et al., 1993; Thiemann et al., 2001; Vrugt et al., 2003b; Vrugt et al.,
346 2008) and may be considered a standard benchmark for parameter estimation of hydrological
347 models. In this regard, three years (i.e., 1953-1955) of hydrologic data (i.e., mean areal
348 precipitation [mm/d], potential evapotranspiration [mm/d], and streamflow [m³/s]) were used.
349 The first two years of data were used for model calibration, while the third year served as a
350 validation dataset for assessing predictive capability. We used the simulation model HYMOD in
351 this catchment to predict streamflow at a single location in the Leaf River channel network. The
352 HYMOD model is a relatively simple rainfall excess model (Moore, 1985) connected with a
353 series of linear reservoirs. HYMOD requires estimation of five parameters and these are listed in
354 Table 1 along with their prior range.

355 **[Table 1 goes here.]**

356 The second case study is the 1,131 km² Hornad River catchment located in Slovakia. The
357 observations for this catchment were collected from 1991 to 2000, and the first five years (i.e.,
358 1991 to 1995) were used for calibration and the remaining data (i.e., 1996 to 2000) was used for
359 validation. We used the simulation model WetSpa in this catchment to predict streamflow at a
360 single location in the Hornad River channel network. Unlike HYMOD, WetSpa is a grid-based
361 hydrologic model that simulates water and energy transfer between soil, plants and the
362 atmosphere. WetSpa can be configured to run in semi-distributed or fully distributed mode of
363 which the former was chosen for this study. According to the previous applications of WetSpa
364 model to Hornad catchment (Bahremand et al., 2007; Liu et al., 2003; Shafii and Smedt, 2009),
365 and as shown in Table 2, 11 WetSpa parameters were targeted for calibration.

366 **[Table 2 goes here.]**

367 The multi-criteria formulation used in this paper was created by splitting a single time series
368 of responses (i.e., discharges) into high- and low-flows. Following Schaepli et al. (2007), high-
369 flows corresponded to time steps in which the hydrograph was rising, and low-flows were
370 defined based on the recession part of hydrograph. Separate Nash-Sutcliffe values (or formal
371 likelihood values, in the case of MCMC sampling) were then calculated for each flow regime,
372 yielding a bi-criteria calibration problem.

373 The computational overhead required for GLUE and DREAM are both dominated by the
374 simulation model run time and as such, for the same number of model simulations completed,
375 GLUE and DREAM require approximately the same computation time. The simulation model
376 run time for HYMOD and WetSpa are 0.65 and 2.25 seconds, respectively, on a PC with 3-GHz
377 Intel processor.

378 **3 Results**

379 For each of the case studies, the DREAM sampler was first applied to establish a converged
380 chain of samples, and the non-converged DREAM and GLUE were then applied. Note that, as
381 mentioned earlier, we used an AR-based Bayesian formulation in this paper. Transformation
382 and/or scaling of parameters is an important factor that can affect the difficulty of parameter
383 estimation (Bates and Watts, 1981; Johnston and Pilgrim, 1976; Kuczera, 1983) and the
384 convergence behaviour of MCMC samplers (Hills and Smith, 1992). For the HYMOD Leaf
385 River and WetSpa Hornad River case studies, a series of preliminary numerical experiments
386 were performed to explore alternative parameter transformations within the DREAM sampler.
387 These experiments indicated that the most suitable transformation was to logarithmically
388 transform HYMOD and WetSpa model parameters and use un-transformed auto-regressive
389 parameters. It should also be noted that, in the formal Bayesian approach, discharges were also
390 transformed logarithmically to stabilize the error variance.

391 **3.1 HYMOD**

392 When applied to the HYMOD Leaf River case study, the DREAM sampler converged after
393 approximately 143000 simulations. The convergence of MCMC sampler was checked using the
394 Gelman-Rubin convergence metric, which was also cross-checked to verify residuals normality
395 (via inspection of a QQ-plot) and non-correlation (via inspection of the auto-correlation
396 function). Furthermore, 1000 out of the last 10000 post-convergence samples were taken from
397 the DREAM chain and used to derive baseline posterior parameter distributions. For the non-
398 converged DREAM approach, a new trial of DREAM was considered up to 10000 simulations of
399 which the last 1000 samples were used to derive corresponding posterior distributions. The

400 GLUE method was applied using the generalized likelihood function Eqs (1-2) considering two
401 scenarios, (i) a budget of 10000 simulations called ‘GLUE Low-budget’, and (ii) identical
402 computation budget to DREAM (i.e., 143000 simulations in HYMOD case study) and is called
403 ‘GLUE Full-budget’.

404 Figure 2 illustrates the posterior parameter information derived by the various calibration
405 methods when applied to the HYMOD Leaf River case study. As observed in Figure 2 the
406 posterior parameter ranges varied across methods, especially with respect to parameters R_S and
407 R_Q . Most of the ranges given by non-converged DREAM were wider than those given by
408 converged DREAM. The difference between the location of posterior solutions derived from
409 Bayesian inference and GLUE is not surprising, and can be explained by the fact that different
410 likelihood functions have been used in these methods. However, comparison between these
411 posterior ranges indicates that incorporating two additional error parameters (i.e., higher
412 complexity in comparison to informal formulation) resulted in a higher level of identifiability,
413 especially for parameters R_S and R_Q .

414 **[Figure 2 goes here.]**

415

416 Figure 3 illustrates the Nash-Sutcliffe (NS) values of the HYMOD Leaf River case study for
417 calibration (upper panel) and validation (lower panel) period, demonstrating the results of
418 DREAM (light points) versus non-converged DREAM and GLUE (dark points) along low and
419 full computational budget. Conversion of DREAM likelihood values into equivalent NS values
420 was non-trivial because the fitted error series should also be accounted for. Proper conversion
421 into equivalent NS values must consider additional elements of the revised Bayesian

422 formulation, namely, the two extra auto-regressive parameters (i.e., ρ_L and ρ_H) and the AR-
423 based residuals term (δ_t). Thus, for a given parameter vector $\boldsymbol{\phi}_i$ containing a model parameter
424 set $\boldsymbol{\theta}_i$ and corresponding $\rho_{L,i}$ and $\rho_{H,i}$ auto-regressive parameters, the corresponding error
425 variances were sampled to generate 100 different time series of error realizations. These errors
426 were then combined with simulated discharges and auto-regressive terms to yield 100 different
427 NS values for parameter vector $\boldsymbol{\phi}_i$. The average of these NS values was then used as the
428 equivalent NS value converted from the original DREAM likelihood value. Repeating this
429 process for all parameter vectors contained in the DREAM posterior samples yielded the
430 equivalent NS values plotted in Figure 3 for calibration and validation period.

431 **[Figure 3 goes here.]**

432 As shown in the calibration part in Figure 3 (upper panel), the results obtained from DREAM
433 were superior (based on NS values) to those given by other methods, and there was some overlap
434 between the posterior sets of solutions given by converged and non-converged DREAM sampler.
435 Note that we sometimes call these sets of solutions ‘posterior clouds’, as they look like a cloud in
436 NS space. In the validation part of Figure 3 (lower panel), the non-converged DREAM posterior
437 cloud very closely resembles the DREAM posterior cloud. This is a good indication that much
438 of the high-density areas of the parameter space were explored prior to the DREAM sampler
439 satisfying the Gelman-Rubin convergence criteria.

440 The results of GLUE in Figure 3 also indicate that regardless of the computation budget
441 considered, the samples were located in fairly identical space in NS space (but with different
442 densities) both in calibration and validation period. However, GLUE with full computational

443 budget performed slightly better considering extreme NS values of GLUE in Figure 3. The
444 GLUE results in calibration period showed that 8% of behavioural samples resulted in negative
445 NS values for low-flows, but since their NS values for high-flows were high, they could rank in
446 the top 10% of all GLUE samples. It should be pointed out that, similar to previous studies
447 (Balin-Talamba, 2004; Lamb et al., 1998), the threshold for classifying solutions as behavioural
448 utilized the formulations in Eqs. (1-2), and did not take into consideration the condition of
449 positive NS values. This explains why there are some solutions with negative low-flows NS
450 values among posterior samples.

451 Figure 3 also shows that GLUE yielded good performance in terms of matching the
452 simulations with observation in validation low-flows, but not as good in high-flows compared to
453 DREAM sampler. In contrast, the posterior cloud generated by DREAM in validation period
454 (Figure 3 lower panel) emphasized matching high-flows (i.e., points clustered in the 0.8 to 1.0
455 range for NS_{high}) at the expense of matching low-flows (i.e., points clustered around $NS_{low} = 0.5$).

456 Ideally, all posterior samples would generate positive NS values in validation period for low-
457 and high-flows. The vertical dashed lines in Figure 3 (lower panel) separates the region with
458 positive NS values for low-flows, and thus, the ideal region would be the right half of the scatter
459 plots. It is observed that all posterior samples from DREAM and all but one of the non-
460 converged DREAM posterior samples were located in this ideal region. However, almost 40% of
461 posterior GLUE (full-budget) samples generated negative validation period NS values for ‘low-
462 flows’. It should be pointed out that almost 92% of these samples had resulted in positive NS
463 values both for low- and high-flows in calibration period.

464 Figure 4 (left panels) illustrates the tradeoff between reliability and sharpness measures for
465 the HYMOD Leaf River case study (only in validation period) for the various methods that were

466 considered (i.e., DREAM, non-converged DREAM, and GLUE with low and full computational
467 budget). The reliability and sharpness values were calculated based on 95% prediction intervals
468 on the corresponding posterior PDFs of simulated discharges. The reliability was calculated
469 based the percentage of coverage of observations by prediction bounds, whereas sharpness was
470 based on the amount of reduction in discharge ranges through comparison with the range of
471 model simulations using prior parameter ranges. In order to define such prior intervals, 100000
472 Latin hypercube samples were taken from prior parameter ranges which were used in HYMOD
473 to generate 100000 discharge hydrographs. The minimum and maximum of discharges at each
474 time steps were then identified to serve as prior discharge ranges.

475 **[Figure 4 goes here.]**

476 The HYMOD results in Figure 4 (left panel) show that the converged DREAM sampler and
477 ‘GLUE Full-budget’ cannot dominate each other with respect to both reliability and sharpness.
478 Compared to ‘GLUE Full-budget’, the converged DREAM resulted in improved sharpness both
479 for low- and high-flows. In terms of reliability, as the goal was to generate 95% prediction
480 intervals, both methods came fairly close to this goal given that reliabilities in validation period
481 ranged from 93% to 97%. Comparison between non-converged DREAM and ‘GLUE Low-
482 budget’ shows that neither of these two methods is superior to the other one with respect to both
483 reliability and sharpness. The reliabilities of these two methods were close to 95%. The
484 sharpness of non-converged DREAM was larger than ‘GLUE Low-budget’ in low-flows, and
485 approximately the same in high-flows.

486 Figure 5 contains Bayesian p-values for both the calibration and validation periods of the
487 HYMOD Leaf River case study for non-converged and converged DREAM approaches. Note

488 that the p-values were derived using the entire set of posterior solutions. Figure 5 shows that
489 even though the p-value results for the converged and non-converged DREAM sampler were
490 different during the calibration period, the results in validation period, however, were fairly
491 similar. Also, both methods yielded underestimation of predictive uncertainty with respect to
492 low-flows in validation period. This might be due to the fact that we used standard Bayesian
493 formulation without disaggregation of different sources of uncertainty, which will be discussed
494 later in the discussion section.

495 **[Figure 5 goes here.]**

496 Figure 6 illustrates the prediction bounds given by the posterior simulations of the considered
497 calibration techniques for the validation period in HYMOD case study. The bounds shown in
498 Figure 6 are derived in a manner similar to those given for posterior parameters of Figure 2 and
499 are assumed to represent 95% prediction intervals. As shown in Figure 6, the converged
500 DREAM sampler reliably covers the validation dataset. Prediction bounds of the non-converged
501 DREAM sampler resemble those generated from the converged DREAM sampler but at the cost
502 of larger width and larger peak flow values. Figure 6 also shows that the prediction bounds
503 associated with ‘GLUE Full-budget’ are larger than those derived with ‘GLUE Low-budget’, but
504 covered the observations better.

505 **[Figure 6 goes here.]**

506 Across the various comparative measures that were evaluated in the context of the HYMOD
507 Leaf River case study, we observed that the formal Bayesian method (both converged and non-
508 converged MCMC sampling) turned out to be more appropriate than informal GLUE strategy in

509 calibration period. Once the validation period was used to evaluate the methods, the formal
510 Bayesian inference (given the formulation of this paper) resulted in a level of underestimation of
511 predictive uncertainty, which would be probably solved through more complex HBS systems, as
512 elaborated in discussions section. On the other hand, the GLUE methodology was only
513 successful in partially meeting the predictive criteria in validation period. The WetSpa Hornad
514 River real case study (Section 3.2) investigates whether these findings would hold for a more
515 complex hydrological model (involving more uncertain parameters) applied to a different
516 catchment.

517 **3.2 WetSpa**

518 For the WetSpa case study (i.e., application to Hornad River catchment), the DREAM
519 sampler was again configured to use a formal auto-regressive Bayesian inference formulation
520 and the method converged (based on the Gelman-Rubin statistic) after 470,000 simulations. As
521 with the HYMOD studies, 10000 post-convergence DREAM samples were taken to construct the
522 Bayesian posterior distributions. Similar to the previous case, the results of non-converged
523 DREAM were derived based on running DREAM only up to 10000 simulations (independent
524 trial than converged DREAM). GLUE was also applied to the WetSpa case study using low and
525 full computational budget as described in HYMOD Leaf River case study.

526 Figure 7 contains normalized posterior ranges of the WetSpa model parameters generated by
527 the various calibration methods. The first result noted in Figure 7 is that some parameters were
528 deemed non-identifiable (i.e., K_S , K_{GI} , and K_{RD}) by the converged DREAM sampler, as indicated
529 by 95% posterior intervals covering almost the entire prior range. When informal likelihood
530 functions were used (i.e., GLUE), most of parameters appeared to be poorly-identifiable.

531 However, it should be noted that the difference between the location of posterior parameter
532 ranges and identifiability levels obtained by formal and informal methods would be explained by
533 the difference in the likelihood functions used in these methods. It is also observed in Figure 7
534 that the posterior parameter ranges derived from non-converged DREAM covered those obtained
535 from converged DREAM, and this shows how the sampler located a smaller posterior region
536 after it converged.

537 **[Figure 7 goes here.]**

538 Figure 8 illustrates the Nash-Sutcliffe values for calibration (upper panel) and validation
539 (lower panel) period of the WetSpa Hornad River case study as evaluated by non-converged
540 DREAM and GLUE (dark points), in comparison to those calculated based on the posterior
541 solutions of the converged DREAM sampler (light points). Note that two cases were reported for
542 GLUE, one with low and one with full computational budget. Also note that the axes in lower
543 panel of Figure 8 were centred between ± 1 , the dashed lines showing the origin where both NS
544 values were zero. A number of GLUE solutions were not within this range and were not depicted
545 in Figure 8. The ideal region for a given calibration method to sample from would be the upper
546 right quadrant of validation panel where both low- and high-flow NS values were positive. It is
547 observed in the calibration panel that DREAM yielded the best NS values both for low- and
548 high-flows. Given that non-converged DREAM and DREAM achieve these high NS values, it
549 seems the inclusion of an error term is important to achieve such high performance. The
550 posterior cloud from non-converged DREAM overlaps substantially the converged DREAM
551 posterior cloud, which indicates that the posterior distribution has likely been sampled from well
552 before the Gelman-Rubin statistic indicated convergence. The results of GLUE (low and full

553 computational budgets) also indicate that increasing the number of simulations in GLUE did not
554 result in comparable model performance as DREAM (see distance between the location of
555 posterior clouds). It is also observed in GLUE results (both low and full budgets) that there were
556 a considerable number of points not located in the ideal region, that is, positive NS values for
557 low and high-flows or the upper right quadrant identified by dashed lines, even though they were
558 all behavioural in the calibration period.

559 **[Figure 8 goes here.]**

560 The sharpness and reliability measures for the validation period of the WetSpa Hornad River
561 case study are given in Figure 4 (right panel). These measures were computed in the same
562 manner as those for the HYMOD Leaf River case study. In terms of reliability, as the goal was to
563 generate 95% prediction intervals, all methods came fairly close to this goal for high flows given
564 that reliabilities in validation period ranged from 94% to 98%. The same is true for validation
565 period low flows except that ‘GLUE Low-budget’ results have a slightly lower reliability of
566 88%. Comparing converged DREAM with ‘GLUE Full-budget’, it is observed that DREAM
567 results dominate GLUE in both low-flows and high flows (i.e., larger reliability and larger
568 sharpness). In other words, DREAM generates tighter 95% prediction intervals and
569 simultaneously improves reliability. Similarly, non-converged DREAM dominates ‘GLUE Low-
570 budget’ results in high flows and practically dominates ‘GLUE Low-budget’ results in low flows
571 (very similar reliabilities but significantly improved sharpness for DREAM).

572 Figure 9 compares the Bayesian p-value QQ plots for non-converged and converged
573 DREAM sampling for the calibration (upper panel) and validation periods (lower panel) of the
574 WetSpa Hornad River case study. As implied by the sigmoid shapes of their respective p-value

575 curves, both DREAM samplers (i.e., converged and non-converged) exhibited systematic under-
576 estimation of uncertainty for low-flows in validation period, even though the results of
577 converged DREAM in calibration period were promising both for low-flows and high-flows.
578 This finding is similar to results in Thyer et al. (2009) and the previous HYMOD case study. The
579 under-estimation of only low-flow uncertainty by the converged DREAM procedure can be
580 considered as indication of model structural error. This suggests that improving the low-flow
581 modules in WetSpa may be a worthwhile enterprise. Such insight highlights the usefulness of
582 multi-criteria Bayesian p-value separation as a post-diagnostic measure for detecting model
583 structural deficiencies. However, it is also possible that the above-mentioned issue may be due to
584 mis-specification of likelihood function.

585 **[Figure 9 goes here.]**

586 Figure 10 illustrates the prediction bounds given by the posterior simulations of the
587 considered calibration techniques for one year (i.e., 1999) of the 5-year validation period
588 (whereas Figure 4 reliability and sharpness values summarize prediction bounds over the entire
589 5-year period). The bounds shown in Figure 10 were derived in a manner similar to those given
590 for posterior parameters of Figure 7 and are assumed to represent 95% prediction intervals. As
591 shown in Figure 10, the converged DREAM sampler reliably covered the validation dataset even
592 though the Bayesian p-value analysis indicated that the results were not perfect with respect to
593 low-flows. Prediction bounds of the non-converged DREAM sampler resemble those generated
594 from the converged DREAM sampler but at the cost of larger width and larger peak flow values.
595 The prediction bounds associated with ‘GLUE Full-budget’ are larger than those derived with
596 ‘GLUE Low-budget’, but covered the observations better.

597

[Figure 10 goes here.]

598 Across all comparative measures, the results of the WetSpa case study suggest the following
599 conclusions: (1) the formal Bayesian inference through the standard formulation of this paper
600 using converged DREAM yielded good results with respect to almost all predictive measures,
601 except for p-values of low-flows in validation period; (2) the non-converged DREAM sampler
602 yielded results that were nearly universally consistent with the converged DREAM sampler
603 while requiring a fraction (i.e., 2%) of the computational budget; and (3) considering the
604 predictive measures addressed in this study, GLUE did not meet all measures as satisfactorily as
605 formal DREAM methodology, even when the full computational budget was considered.

606 **4 Discussion**

607 The DREAM results suggest that the Gelman-Rubin convergence criterion is too stringent
608 since non-converged DREAM results closely approximates converged DREAM results and yet
609 requires a fraction of the computational budget. It may also be possible to further improve the
610 results of the non-converged DREAM sampler (i.e., make it more closely approximate the
611 converged DREAM results) by filtering out obviously low quality solutions for the calibration
612 period (e.g., those with NS values smaller than 0.5 in upper left panels of Figures 3 and 8). Also,
613 one might think of applying alternative convergence measures. A potential hydrology-based
614 convergence metric can be the reproduction of hydrological signatures that represent the overall
615 hydrologic behaviour of the catchment (Gupta et al., 2008; Yilmaz et al., 2008). Future research
616 should explore these and other alternative convergence measures in a multi-criteria context.

617 Comparison between formal and informal methods could also be viewed from the standpoint
618 of aleatory and epistemic uncertainties, which was also elaborated in the introduction section of
619 this paper. The errors in the case studies of this paper are assumed to be aleatory (especially in
620 Bayesian inference methodology), even though in reality they could be a mixture of both
621 aleatory and epistemic uncertainties. The results reveal that validation period performance
622 measures are generally poorer compared to calibration period which is expected to be caused by
623 epistemic errors (Beven et al., 2011). Thus, in the presence of epistemic errors, neither the
624 standard Bayesian formulation nor the informal methods (such as GLUE) would be perfectly
625 reliable in prediction mode. There are improved informal and formal approaches for case studies
626 where epistemic errors are thought to be significant, e.g., the use of hierarchical Bayesian
627 structures (e.g., Huard, 2008; Kuczera et al., 2006; Moradkhani et al., 2005; Renard et al., 2010;
628 Wei et al., 2010), or the concept of ‘limits of acceptability’ used for identifying behavioural
629 models in GLUE (Blazkova and Beven, 2009; Liu et al., 2009). Comparison between these two
630 more advanced formal and informal uncertainty analysis methods is an interesting future
631 research avenue.

632 **5 Concluding Remarks**

633 This paper evaluates the applicability of formal (Bayesian inference) and informal (GLUE)
634 multi-criteria methods to uncertainty-based calibration in hydrological modelling. Bayesian
635 inference is implemented through DREAM sampling based on a multi-criteria formulation. The
636 results of non-converged DREAM are also evaluated. The results are compared with those
637 obtained from two scenarios for GLUE, using a restricted computational budget and the full
638 computational budget equivalent to the budget required for DREAM sampler to converge. The

639 various methods are applied to two cases involving the 5-parameter HYMOD model and the 11-
640 parameter WetSpa model. Results demonstrate that there can be considerable differences in
641 prediction intervals generated by formal and informal strategies for uncertainty-based multi-
642 criteria calibration. Future uncertainty-based calibration studies for simulation models with a
643 large number of parameters should be aware of the potential considerable difference between the
644 results of formal and informal strategies.

645 Results also demonstrate that it is advisable to consider multiple comparative measures,
646 including traditional metrics like the Nash-Sutcliffe efficiency, when comparing alternative
647 calibration strategies. Furthermore, it is observed that the choice of using the validation period or
648 the calibration period for selected comparative measures would influence the analysis and as
649 such it is recommended that future uncertainty-based calibration method comparison studies
650 should include and largely focus on comparative performance assessment for the validation
651 period.

652 In general, the Bayesian inference methodology performs well (in comparison with other
653 methods) along all comparative measures except for low-flows in validation period considering
654 the same computational budget, e.g., DREAM validation period prediction intervals are
655 simultaneously tighter and more reliable than corresponding GLUE intervals. In case of limited
656 computational budget (i.e., only 10000 simulations in this paper), non-converged MCMC
657 sampling using DREAM proves to be fairly consistent with formal Bayesian inference. This
658 indicates the potential value of utilizing formal MCMC sampling results before convergence as a
659 promising alternative to informal methods such as GLUE.

660 The results obtained through application of Bayesian inference to the two cases of this paper
661 indicated under-estimation of predictive uncertainty for low-flows in the validation period. We

662 applied a standard Bayesian formulation which lumps all uncertainties into a single additive error
663 term. More recently, Renard et al. (2010; 2011) showed that consideration of rainfall and model
664 structural uncertainties outside of the error term used in Bayesian formulation yielded more
665 reliable estimation of the predictive uncertainty for all runoff ranges, as opposed to the typical
666 Bayesian formulation in our paper. Application of hierarchical Bayesian structures to the case
667 studies of this paper is currently being investigated.

668 There are many ways to formulate and conduct GLUE analyses, and to some extent DREAM
669 calibration experiments. Our experiments require a number of subjective decisions and as such
670 our results are conditional on these decisions. However, we believe that the subjective decisions
671 we make are consistent with the decisions others have made in the literature. For example,
672 although it is possible to apply GLUE using a formal likelihood function, the literature suggests
673 that is relatively uncommon and thus we do not examine this. We used GLUE with an informal
674 generalized likelihood function in this study because the objective of the study was to assess its
675 performance as an informal method. It may be possible that applying informal methods such as
676 GLUE using formal likelihood functions would improve their performance, but this is not the
677 focus of the present study. Future comparative studies systematically varying such subjective
678 decisions would be valuable.

679 **6 Acknowledgement**

680 The authors would like to acknowledge Dr. Jasper A. Vrugt for providing the code of his
681 DREAM algorithm used in this study and the four anonymous reviewers for their helpful
682 comments that have improved the paper.

683

684 **References**

685 Bahremand, A. et al., 2007. WetSpa Model Application for Assessing Reforestation Impacts on Floods in
686 Margecany–Hornad Watershed, Slovakia. *Water Resources Management*, 21(8): 1373-1391.

687 Balin-Talamba, D., 2004. Hydrological behaviour through experimental and modeling approaches;
688 Application to the Haute-Mentue catchment, PhD Thesis, Swiss Federal School of Technology of
689 Lausanne.

690 Balin-Talamba, D., Parent, E., Musy, A., 2010. Bayesian multiresponse calibration of TOPMODEL:
691 Application to the Haute-Mentue catchment, Switzerland. *Water Resources Research*, 46:
692 W08524.

693 Bates, B.C., Campbell, E.P., 2001. A Markov Chain Monte Carlo Scheme for parameter estimation and
694 inference in conceptual rainfall-runoff modeling. *Water Resources Research*, 37(4): 937-947.

695 Bates, D.M., Watts, D.G., 1981. Parameter transformations for improved approximate confidence
696 regions in nonlinear least squares. *The Annals of Statistics*, 9(6): 1152-1167.

697 Bates, D.M., Watts, D.G., 1988. *Nonlinear regression analysis and its applications*. Wiley, NY.

698 Beven, K., Binley, A., 1992. The future of distributed models: model calibration and uncertainty
699 prediction. *Hydrological Processes*, 6(3): 279-298.

700 Beven, K., Freer, J., 2001. Equifinality, data assimilation, and uncertainty estimation in mechanistic
701 modelling of complex environmental systems using the GLUE methodology. *J. Hydrology*,
702 249((1-4)): 11-29.

703 Beven, K., Smith, P.J., Westerberg, I., Freer, J., 2012. Comment on “Pursuing the method of multiple
704 working hypotheses for hydrological modeling” by P. Clark et al. *Water Resources Research*, 48:
705 W11801.

706 Beven, K., Smith, P.J., Wood, A., 2011. On the colour and spin of epistemic error (and what we might do
707 about it). *Hydrol. Earth Syst. Sci.*, 15: 3123-3133.

708 Beven, K.J., 2009. Comment on "Equifinality of formal (DREAM) and informal (GLUE) Bayesian
709 approaches in hydrologic modeling; by Jasper A. Vrugt, Cajo J. F. ter Braak, Hoshin V. Gupta and
710 Bruce A. Robinson. *Stoch Environ Res Risk Assess*, 23: 1059-1060.

711 Beven, K.J., Smith, P.J., Freer, J., 2008. So just why would a modeller choose to be incoherent? . *J.*
712 *Hydrology*, 354: 15-32.

713 Binley, A., Beven, K., 1991. Physically-based modelling of catchment hydrology: a likelihood approach to
714 reducing predictive uncertainty. In: Rycroft, D.G.F.a.M.J. (Ed.), *Computer Modelling in the*
715 *Environmental Sciences, The Institute of Mathematics and its Applications Conference Series*.
716 Clarendon Press, Oxford, pp. 75–88.

717 Blasone, R.-S., Madsen, H., Rosbjerg, D., 2008a. Uncertainty assessment of integrated distributed
718 hydrological models using GLUE with Markov chain Monte Carlo sampling. *J. Hydrology*, 353:
719 18– 32.

720 Blasone, R.S. et al., 2008b. Generalized likelihood uncertainty estimation (GLUE) using adaptive Markov
721 Chain Monte Carlo sampling. *Adv. Water Res.*, 31: 630-648.

722 Blazkova, S., Beven, K., 2009. A limits of acceptability approach to model evaluation and uncertainty
723 estimation in flood frequency estimation by continuous simulation: Skalka catchment, Czech
724 Republic. *Water Resources Research*, 45: W00B16.

725 Box, G.E.P., Tiao, G.C., 1973. *Bayesian Inference in Statistical Analysis*. Addison-Wesley, Boston, MA.

726 Boyle, D.P., 2000. Multicriteria calibration of hydrological models. Ph.D. Dissertation Thesis, Univ. of
727 Ariz.

728 Clark, M.P., Kavetski, D., Fenicia, F., 2012. Reply to comment by K. Beven et al. on “Pursuing the method
729 of multiple working hypotheses for hydrological modeling”. *Water Resources Research*, 48:
730 W11802.

731 Freni, G., Mannina, G., 2009. Bayesian approach for uncertainty quantification in water quality
732 modelling: The influence of prior distribution. *J. Hydrology*, 392: 31–39.

733 Gelman, A., Carlin, J.B., Stern, H.S., Rubin, D.B., 2004. *Bayesian Data Analysis*. Chapman & Hall/CRC, Boca
734 Raton, Florida.

735 Gupta, H.V., Sorooshian, S., Yapo, P.O., 1998. Toward improved calibration of hydrologic models:
736 Multiple and noncommensurable measures of information. *Water Resources Research*, 34(4):
737 751-764.

738 Gupta, H.V., Wagener, T., Liu, Y., 2008. Reconciling theory with observations: Elements of a diagnostic
739 approach to model evaluation. *Hydrological Processes*.

740 Hamilton, S., 2007. Just say NO to equifinality. *Hydrological Processes*, 21(14): 1979-1980.

741 Hills, S.E., Smith, A.F.M., 1992. Parameterization issues in Bayesian inference. In: J. M. Bernardo, J.B., A.
742 P. Dawid and A. F. M. Smith (Ed.), *Bayesian Statistics 4*. Oxford University Press, pp. 227-46.

743 Hong, B., Strawderman, R.L., Swaney, D.P., Weinstein, D.A., 2005. Bayesian estimation of input
744 parameters of a nitrogen cycle model applied to a forested reference watershed, Hubbard Brook
745 Watershed Six. *Water Resources Research*, 41.

746 Huard, D., and A. Mailhot, 2008. Calibration of hydrological model GR2M using Bayesian uncertainty
747 analysis. *Water Resources Research*, 44: W02424.

748 Jeremiah, E., Sisson, S., Marshall, L., Mehrotra, R., Sharma, A., 2011. Bayesian calibration and
749 uncertainty analysis of hydrological models: A comparison of adaptive Metropolis and
750 sequential Monte Carlo samplers. *Water Resources Management*, 47(W07547).

751 Jin, X., Xu, C.-Y., Zhang, Q., Singh, V.P., 2010. Parameter and modeling uncertainty simulated by GLUE
752 and a formal Bayesian method for a conceptual hydrological model. *Journal of Hydrology*, 383:
753 147-155.

754 Johnston, P.R., Pilgrim, D.H., 1976. Parameter optimization for Watershed models. *Water Resources
755 Research*, 12(3): 477-486.

756 Kavetski, D., Franks, S.W., Kuczera, G., 2002. Confronting Input Uncertainty in Environmental Modeling.
757 In: Duan, Q., Gupta, H.V., Sorooshian, S., Rousseau, A.N., Turcotte, R. (Eds.), *Calibration of
758 Watershed Models*. AGU, Washington, DC, pp. 49-68.

759 Kavetski, D., Kuczera, G., Franks, S.W., 2006. Bayesian analysis of input uncertainty in hydrological
760 modeling: 1. Theory. *Water Resources Research*, 42, W03407, doi: 10.1029/2005WR004368.

- 761 Kuczera, G., 1983. Improved parameter inference in catchment models: 1. Evaluating parameter
762 uncertainty. *Water Resources Research*, 19(5): 1151-1162.
- 763 Kuczera, G., Kavetski, D., Franks, S., Thyer, M., 2006. Towards a Bayesian total error analysis of
764 conceptual rainfall-runoff models: Characterising model error using storm-dependent
765 parameters. *Journal of Hydrology*, 331(3-4): 161–177.
- 766 Kuczera, G., Mroczkowski, M., 1998. Assessment of hydrologic parameter uncertainty and the worth of
767 multiresponse data. *Water Resources Research*, 34(6): 1481–1489.
- 768 Kuczera, G., Parent, E., 1998. Monte Carlo assessment of parameter uncertainty in conceptual
769 catchment models: the Metropolis algorithm. *Journal of Hydrology*, 211(1-4): 69-85.
- 770 Laio, F., Tamea, S., 2007. Verification tools for probabilistic forecasts of continuous hydrological
771 variables. *Hydrol. Earth Syst. Sci.*, 11(4): 1267- 1277.
- 772 Lamb, R., Beven, K.J., Myrabø, S., 1998. Use of spatially distributed water table observations to constrain
773 uncertainty in a rainfall–runoff model. *Adv. Water Res.*, 22(4): 305–317.
- 774 Legates, D.R., McCabe, G.J., 1999. Evaluating the use of “goodness-of-fit” measures in hydrologic and
775 hydro-climatic model evaluation. *Water Resources Research*, 35: 233-241.
- 776 Li, L., Xia, J., Xu, C.-Y., Singh, V.P., 2010. Evaluation of the subjective factors of the GLUE method and
777 comparison with the formal Bayesian method in uncertainty assessment of hydrological models.
778 *Journal of hydrology*, 390: 210-221.
- 779 Liu, Y., Freer, J., Beven, K., Matgen, P., 2009. Towards a limits of acceptability approach to the calibration
780 of hydrological models: Extending observation error. *J. Hydrology*, 367: 93-103.
- 781 Liu, Y.B., Gebremeskel, S., De Smedt, F., Hoffmann, L., Pfister, L., 2003. A diffusive transport approach
782 for flow routing in GIS-based flood modeling. *Journal of Hydrology*, 283(1-4): 91-106.
- 783 Madsen, H., 2000. Automatic calibration of a conceptual rainfall-runoff model using multiple objectives.
784 *J. Hydrology(235)*: 276– 288.
- 785 Mantovan, P., Todini, E., 2006. Hydrological forecasting uncertainty assessment: Incoherence of the
786 GLUE methodology. *J. Hydrology*, 130((1 -2)): 368- 381.
- 787 Marshall, L., Nott, D., Sharma, A., 2007. Towards dynamic catchment modelling: a Bayesian hierarchical
788 mixtures of experts framework. *Hydrological Processes*, 21: 847-861.
- 789 Montanari, A., 2005. Large sample behaviors of the generalized likelihood uncertainty estimation (GLUE)
790 in assessing the uncertainty of rainfall-runoff simulations. *Water Resources Research*, 41:
791 W08406.
- 792 Montanari, A., 2007. What do we mean by "uncertainty"? The need for a consistent wording about
793 uncertainty in hydrology. *HPToday*, 21(6): 841-845.
- 794 Montanari, A., 2011. Interactive comment on “On the colour and spin of epistemic error (and what we
795 might do about it)” by K. Beven et al. *Hydrol. Earth Syst. Sci. Discussion*, 8: C2885–C2891.
- 796 Moore, R.J., 1985. The probability-distributed principle and runoff production at point and basin scales.
797 *Hydrological Sciences*, 30(2): 273–297.
- 798 Moradkhani, H., Sorooshian, S., Gupta, H.V., Houser, P.R., 2005. Dual state-parameter estimation of
799 hydrological models using ensemble Kalman filter. *Adv. Water Resour.*, 28: 135-147.

800 Mroczkowski, M., Raper, G.P., Kuczera, G., 1997. The quest for more powerful validation of conceptual
801 catchment models. *Water Resources Research*, 33: 2325- 2335.

802 Qian, S.S., Stow, C.A., Borsuk, M.E., 2003. On Monte Carlo methods for Bayesian inference. *Ecological*
803 *Modelling*, 159: 269-277.

804 Razavi, S. et al., 2010. Reducing the Computational Cost of Automatic Calibration through Model Pre-
805 Emption. *Water Resources Research*, 46(11): W11523.

806 Refsgaard, J.C., 1997. Parameterisation, calibration and validation of distributed hydrological models. *J.*
807 *Hydrology*, 198: 69-97.

808 Reichert, P., 1997. On the necessity of using imprecise probabilities for modelling environmental
809 systems. *Water Science and Technology*, 36(5): 149-156.

810 Renard, B., Kavetski, D., Kuczera, G., Thyer, M., 2010. Understanding predictive uncertainty in hydrologic
811 modeling: The challenge of identifying input and structural errors. 2010, 46: W05521.

812 Renard, B. et al., 2011. Toward a reliable decomposition of predictive uncertainty in hydrological
813 modeling : Characterizing rainfall errors using conditional simulation. *Water Resour. Res.*, 47:
814 W11516.

815 Romanowicz, R.J., K.J. Beven, Tawn, J., 1994. Evaluation of predictive uncertainty in nonlinear
816 hydrological models using a Bayesian approach. In: Turkman, V.B.a.K.F. (Ed.), *Statistics for the*
817 *Environment 2, Water Related Issues*, pp. 297-315.

818 Schaeffli, B., Talamba, D.B., Musy, A., 2007. Quantifying hydrological modeling errors through a mixture
819 of normal distributions. *J. Hydrology*(332): 303– 315.

820 Seber, G.A., Wild, C.J., 1989. *Nonlinear Regression*. John Wiley and Sons, New York (NY).

821 Shafii, M., Smedt, F.D., 2009. Multi-objective calibration of a distributed hydrological model (WetSpa)
822 using a genetic algorithm. *Hydrology and Earth System Sciences*, 13: 2137-2149.

823 Sorooshian, S., Duan, Q., Gupta, V.K., 1993. Calibration of rainfall-runoff models: Application of global
824 optimization to the Sacramento Soil Moisture accounting model. *Water Resources Research*, 29:
825 1185– 1194.

826 Stedinger, J.R., Vogel, R.M., Lee, S.U., Batchelder, R., 2008. Appraisal of the generalized likelihood
827 uncertainty estimation (GLUE) method. *Water Resources Research*, 44: W00B06.

828 Thiemann, M., Trosset, M., Gupta, H.V., Sorooshian, S., 2001. Bayesian recursive parameter estimation
829 for hydrologic models. *Water Resources Research*, 37(10): 2521-2535.

830 Thyer, M. et al., 2009. Critical evaluation of parameter consistency and predictive uncertainty in
831 hydrological modeling: A case study using Bayesian total error analysis. *Water Resources*
832 *Research*, 45: W00B14.

833 Vrugt, J.A., Gupta, H.V., Bastidas, L.A., Bouten, W., Sorooshian, S., 2003a. Effective and efficient
834 algorithm for multiobjective optimization of hydrologic models. *Water Resources Research*,
835 39(8): 1214, doi:10.1029/2002WR001746.

836 Vrugt, J.A., Gupta, H.V., Bouten, W., Sorooshian, S., 2003b. A Shuffled Complex Evolution Metropolis
837 algorithm for optimization and uncertainty assessment of hydrologic model parameters. *Water*
838 *Resources Research*, 39(8): 1201, doi:10.1029/2002WR001642.

839 Vrugt, J.A. et al., 2009. Accelerating Markov chain Monte Carlo simulation by differential evolution with
840 self-adaptive randomized subspace sampling. *International Journal of Nonlinear Sciences and*
841 *Numerical Simulation*, 10(3): 273-290.

842 Vrugt, J.A., ter Braak, C.J.F., Gupta, H.V., Robinson, B.A., 2008. Equifinality of formal (DREAM) and
843 informal (GLUE) Bayesian approaches in hydrologic modeling? *Stoch Environ. Res. Risk Assess.*,
844 44: 1-16.

845 Wagener, T. et al., 2001. A framework for development and application of hydrological models. *Hydrol.*
846 *Earth Syst. Sci.*, 5(1): 13– 26.

847 Wang, Z.-M., Batelaan, O., De Smedt, F., 1996. A distributed model for water and energy transfer
848 between soil, plants and atmosphere (WetSpa). *Physics and Chemistry of The Earth*, 21(3): 189-
849 193.

850 Wei, W., Clark, J.S., Vose, J.M., 2010. Assimilating multi-source uncertainties of a parsimonious
851 conceptual hydrological model using hierarchical Bayesian modeling. *J. Hydrology*, 394: 436–
852 446.

853 Yadav, M., Wagener, T., Gupta, H., 2007. Regionalization of constraints on expected watershed response
854 behavior for improved predictions in ungauged basins. *Advances in Water Resources*, 30 (8):
855 1756–1774.

856 Yang, J., Reichert, P., Abbaspour, K.C., Yang, H., Xia, J., 2008. Comparing Uncertainty Analysis Techniques
857 for a SWAT Application to the Chaohe Basin in China. *J. Hydrology*, 358: 1-23.

858 Yapo, P.O., Gupta, H.V., Sorooshian, S., 1998. Multi-objective global optimization for hydrologic models.
859 *Journal of Hydrology*, 204(1-4): 83-97.

860 Yilmaz, K.K., Gupta, H.V., Wagener, T., 2008. A process-based diagnostic approach to model evaluation:
861 Application to the NWS distributed hydrologic model. *Water Resources Research*, 44: W09417.

862

863 **APPENDIX – Review of Bayesian Inference Procedure**

864 This appendix provides a summary of the Bayesian formulation used in this paper, and the
865 details can be found in previous studies (Balin-Talamba et al., 2010; Schaeffli et al., 2007). We
866 assume the AR-based formulation as follows:

867
$$Y_i = (Y_i^{sim} | \boldsymbol{\theta}, \mathbf{X}) + \rho \varepsilon_{i-1} + \delta_i \tag{A1}$$

868 where Y_i and Y_i^{sim} are the observed and simulated values for the model response at time step
869 i , $\boldsymbol{\theta}$ is the model parameters vector, \mathbf{X} is the model inputs vector, ρ is the lag-one AR
870 parameter, $\varepsilon_i = (Y_i - Y_i^{sim}(\boldsymbol{\theta}, \mathbf{X}))$ is the residual between observation and model prediction at time
871 step i (and $\varepsilon_0 = 0$), and δ_i is random error term:

872 $\delta_i \sim N(0, \sigma_j^2)$ (A2)

873 with σ_j^2 being the residual variance for response j , here considered unknown and should be
 874 estimated. If we consider J responses, then J parameters (representing error variance for J
 875 responses) need to be estimated in the Bayesian inference methodology. Under the assumption of
 876 multiple and statistically independent responses, the combined statistical likelihood function for
 877 multiple responses is simply the product of the individual likelihood functions:

$$\begin{aligned}
 l_{multiple} &= \prod_{j=1}^J l_j(\boldsymbol{\theta}, \rho, \sigma_j^2, \mathbf{X}) \\
 &= \prod_{j=1}^J \frac{1}{(\sqrt{2\pi})^{t_j} \cdot \sigma_j^{t_j}} \cdot \exp\left(-\frac{\sum_{i=1}^{t_j} \delta_{j,i}^2}{2\sigma_j^2}\right)
 \end{aligned}
 \tag{A3}$$

879 where $\delta_{j,i} = \varepsilon_{j,i} - \rho\varepsilon_{j,i-1}$ for observation set j and time step i (note that $\varepsilon_{j,0} = 0$),
 880 respectively; J is the number of observation sets, and t_j is the number of time steps for each
 881 observation set j . In order to derive the posterior distribution of parameters, a bounded uniform
 882 prior distribution is considered for $\boldsymbol{\theta}$ over prior feasible range, and the prior distribution of error
 883 variance is also considered to be Jeffrey non-informative distribution as follows:

884 $p(\sigma_j^2) \propto 1/\sigma_j^2$ for $0 < \sigma_j^2 < \infty$ (A4)

885 Using such prior distributions enables us to integrate out the error variances, and the
 886 Bayesian formulation results in the joint posterior distributions from which the marginal
 887 distribution of model parameters and error variances can be estimated conditioned on the
 888 observed data \mathbf{Y} . Alternatively, we can use MCMC sampling to directly take samples from the
 889 posterior distributions, all of which are contained in the chain. In MCMC implementations, the
 890 acceptance/rejection criterion ratio (between posterior densities of the new candidate and old
 891 current samples) is used to accept/reject the candidate to be added to the chain. In the multi-

892 criteria Bayesian formulation, let $\sigma_{j,current}^2$ and $\sigma_{j,candidate}^2$ be the error variance of the current and
893 candidate solutions, respectively, which are estimated based on the residuals after running the
894 simulation model. Also assume the quantity $S_j = 0.5 \sum_{i=1}^{t_j} \delta_{j,i}^2$, such that $S_{j,current}$ and $S_{j,candidate}$
895 be the values for the current and the candidate solutions, respectively. The final form of the
896 acceptance/rejection criterion can then be shown as follows:

$$897 \quad \alpha = \prod_{i=1}^J \exp \left[\left(\frac{1}{\sigma_{j,current}^2} + \frac{1}{\sigma_{j,candidate}^2} \right) (S_{j,current} - S_{j,candidate}) \right] \cdot \left(\frac{S_{j,candidate}}{S_{j,current}} \right)^{\frac{t_j}{2}} \quad (A5)$$

898
899
900

901

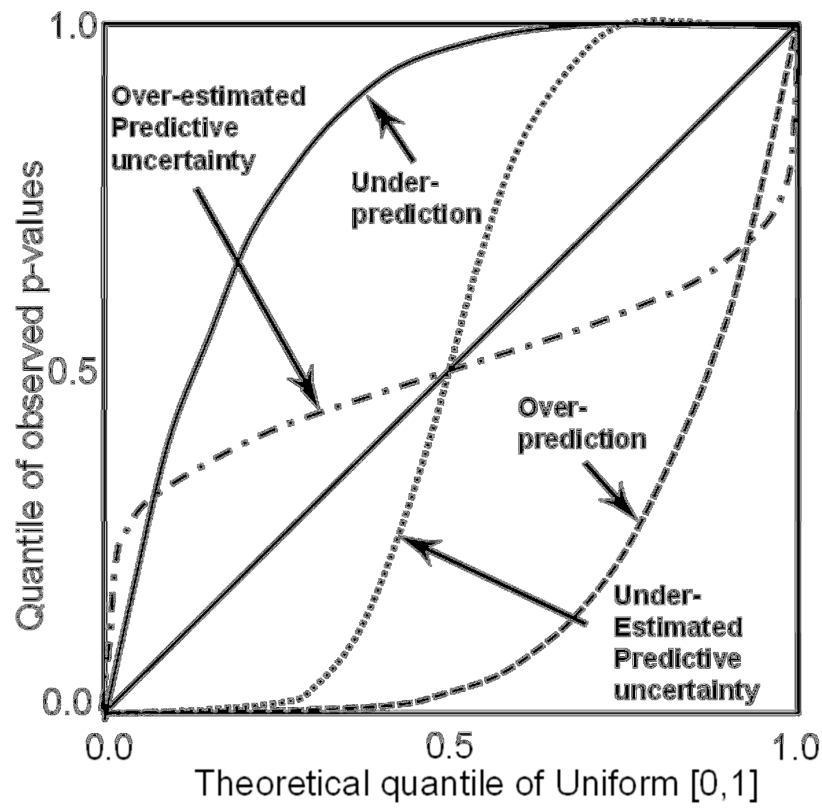
Table 1. HYMOD parameters and their prior range

Parameter	Description	Unit	Prior Range
CMAX	Maximum storage capacity	mm	[1 , 500]
BEXP	Degree of the soil spatial variability moisture capacity	-	[0.1 , 2]
ALPHA	Distributing factor on flow between the two series of reservoirs	-	[0 , 0.1]
RQ	Residence time of the quick reservoirs	d	[0 , 0.1]
RS	Residence time of the slow reservoirs	d	[0.1 , 0.99]

902

Table 2. Parameters of WetSpa simulation model

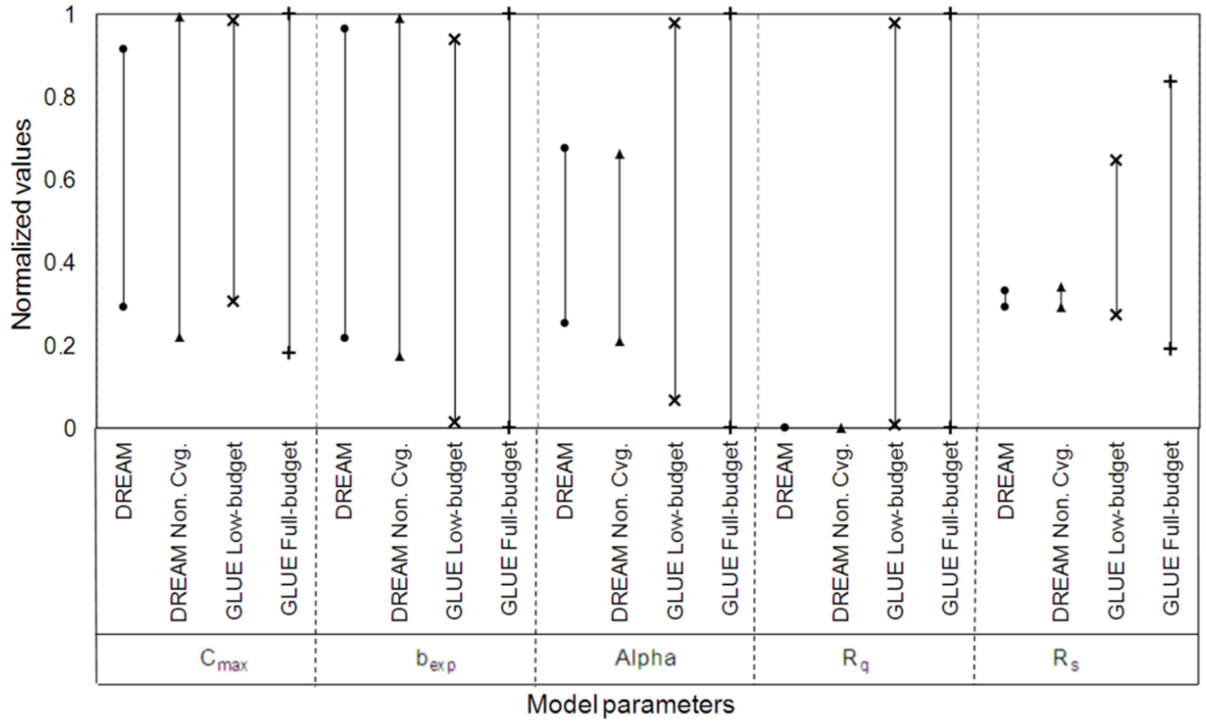
Parameter	Description	Unit	Prior Range
Ki	Interflow scaling factor	-	[0 – 10]
Kg	Groundwater recession coefficient	d ⁻¹	[0 - 0.05]
Ks	Initial soil moisture factor	-	[0 – 2]
Ke	Correction factor for PET	-	[0 – 2]
Kgi	Initial groundwater storage	mm	[0 – 500]
Kgm	Groundwater storage scaling factor	mm	[0 – 2000]
Kt	Base temperature for snowmelt	°C	[-1 – 1]
Ktd	Temperature degree-day coefficient	mm °C ⁻¹ d ⁻¹	[0 – 10]
Krd	Rainfall degree-day coefficient	°C ⁻¹ d ⁻¹	[0 - 0.05]
Km	Surface runoff coefficient	-	[0 – 5]
Kp	Rainfall scaling factor	mm	[0 – 500]



905

906 **Figure 1. Schematic of the predictive QQ plot based on Thyer et al. (2009)**

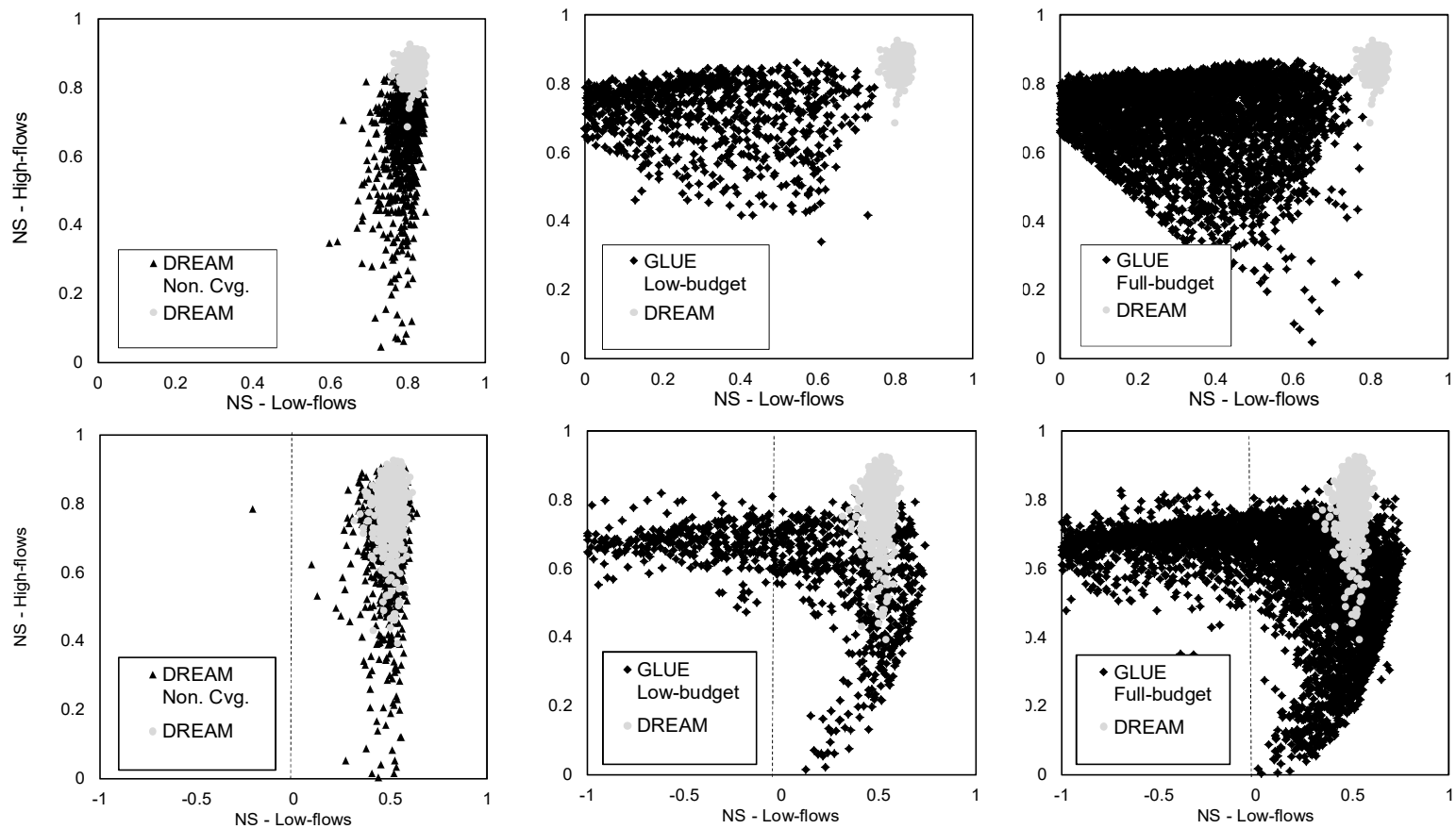
907



908

909 **Figure 2. Posterior ranges of HYMOD parameters for the Leaf River case study; The parameter**
 910 **ranges correspond to 95% posterior intervals for different uncertainty analysis methods.**

911

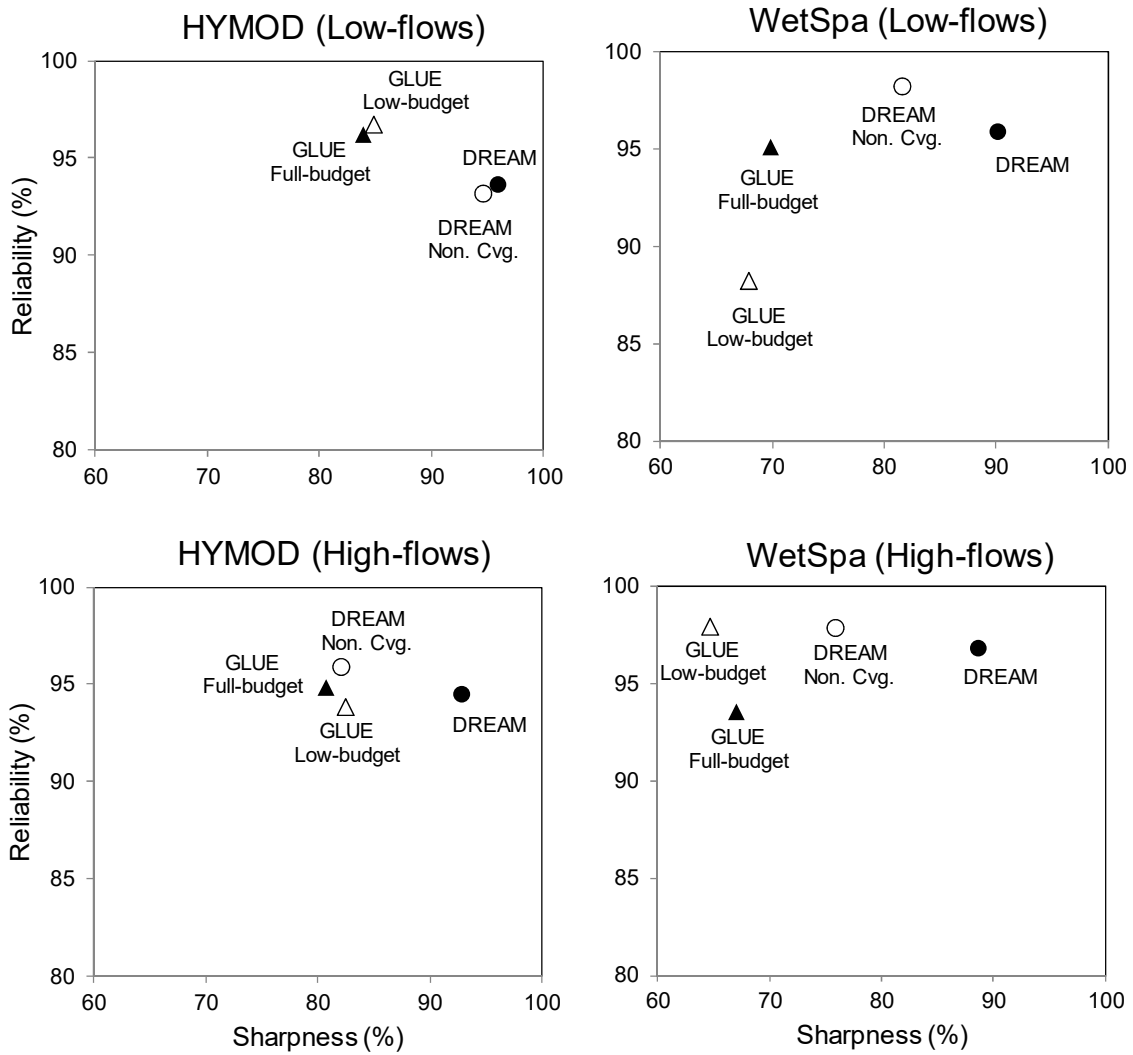


912

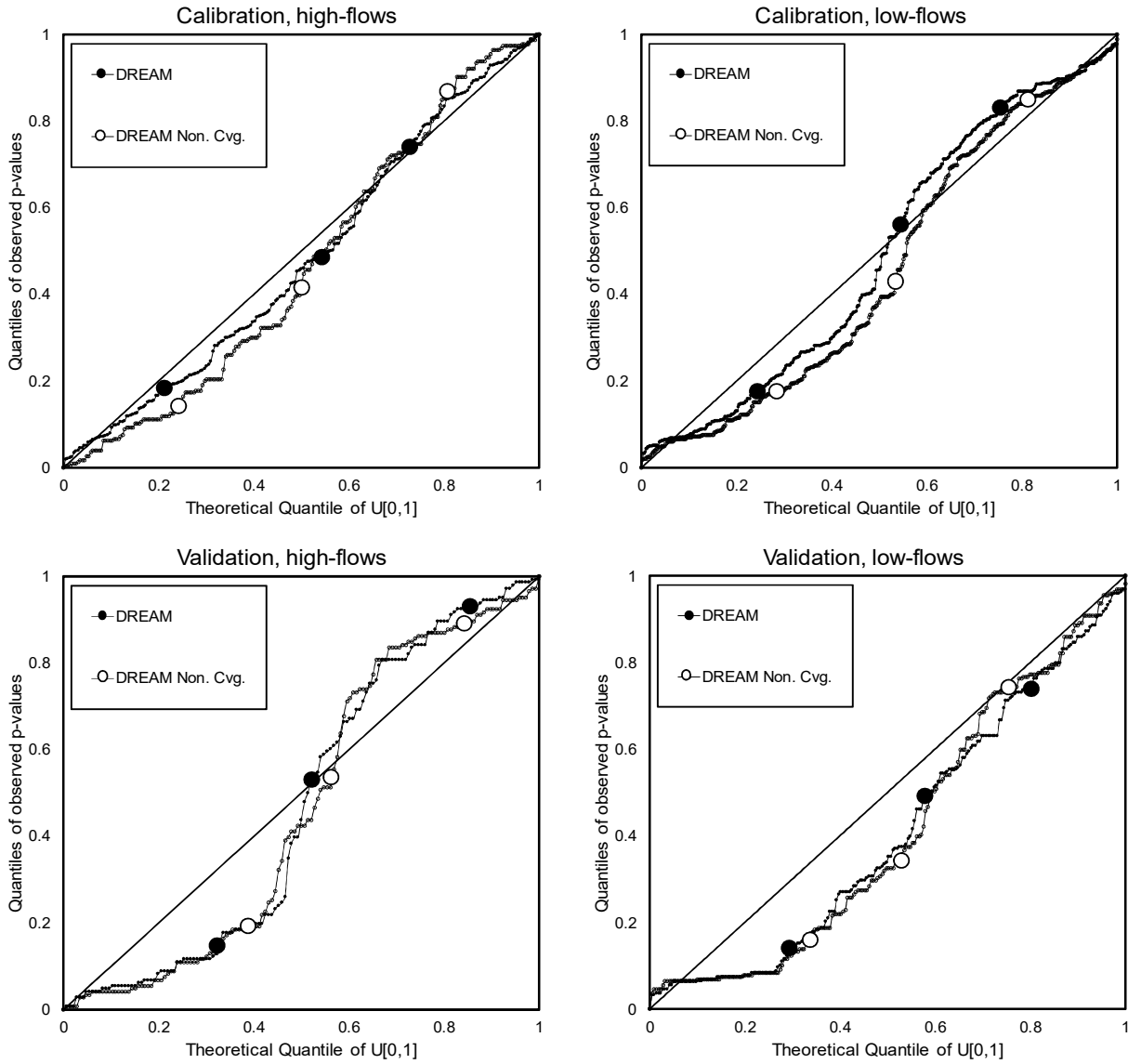
913

914 **Figure 3. NS values of low-flows (horizontal axis) and high-flows (vertical axis) in calibration (upper panels) and validation (lower**
 915 **panels) period for HYMOD case study, derived from DREAM (light points) versus non-converged DREAM and GLUE methods (dark**
 916 **points).**

917



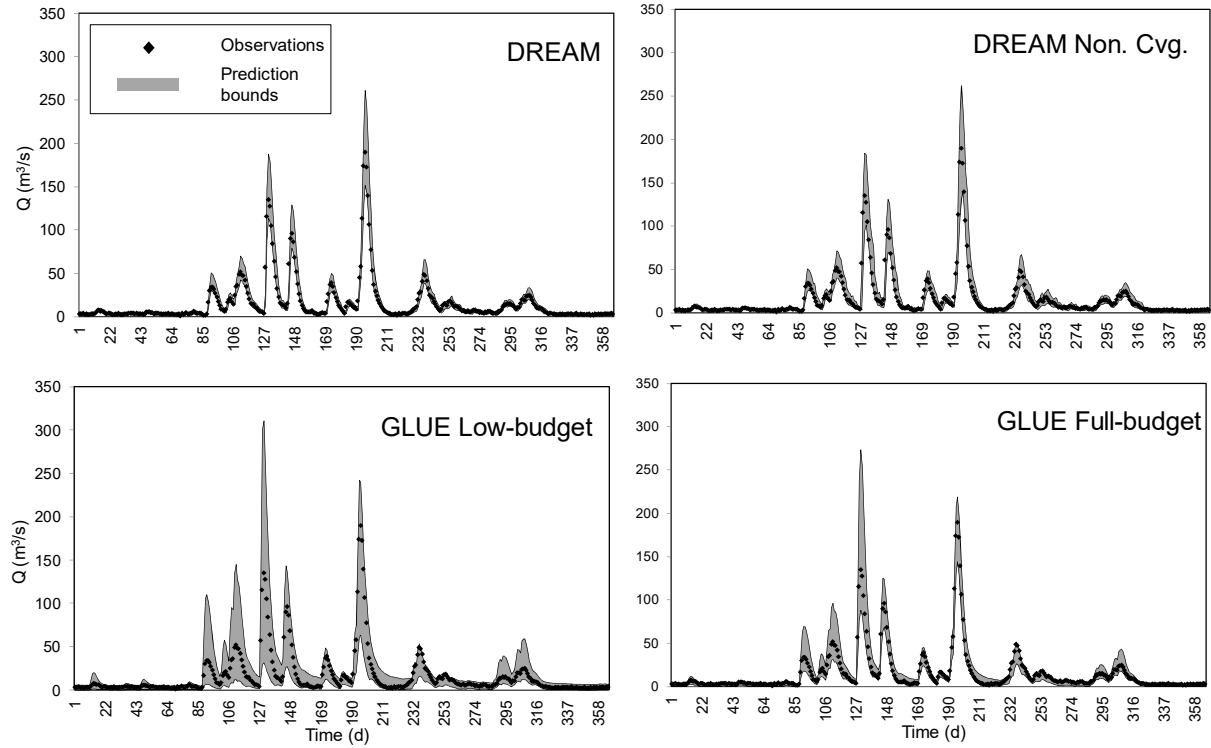
920 **Figure 4. Validation period reliability and sharpness for low-flows (upper panels) and high-flows**
 921 **(lower panels) in application of different techniques (shown in different shapes) to the HYMOD**
 922 **and WetSpa simulation models.**



923

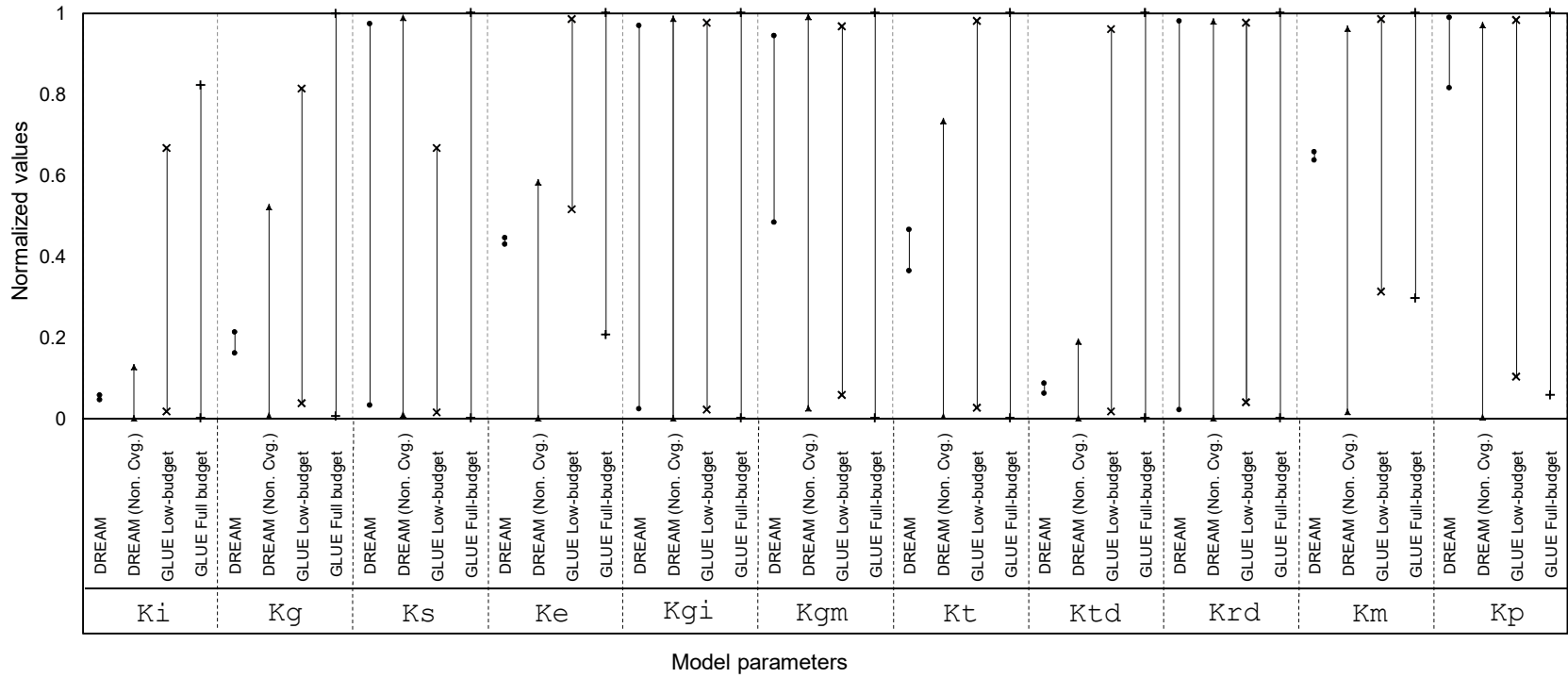
924 **Figure 5. QQ plot of Bayesian p-values for high- and low-flows derived from converged and non-**
 925 **converged DREAM, for calibration (upper) and validation (bottom) periods of the HYMOD Leaf**
 926 **River case study.**

927



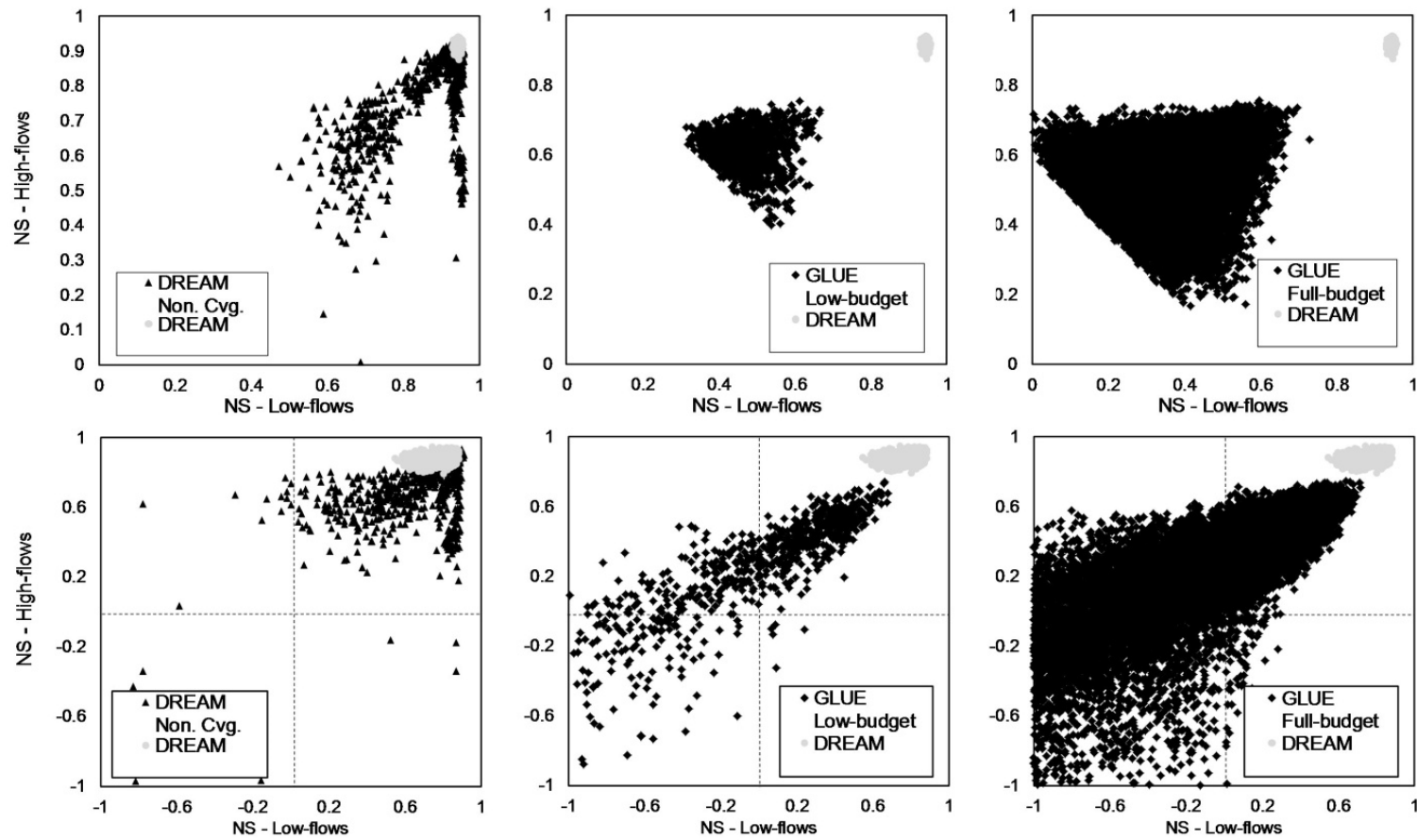
928

929 **Figure 6. Prediction bounds and observations for the validation period in the HYMOD case study.**



930

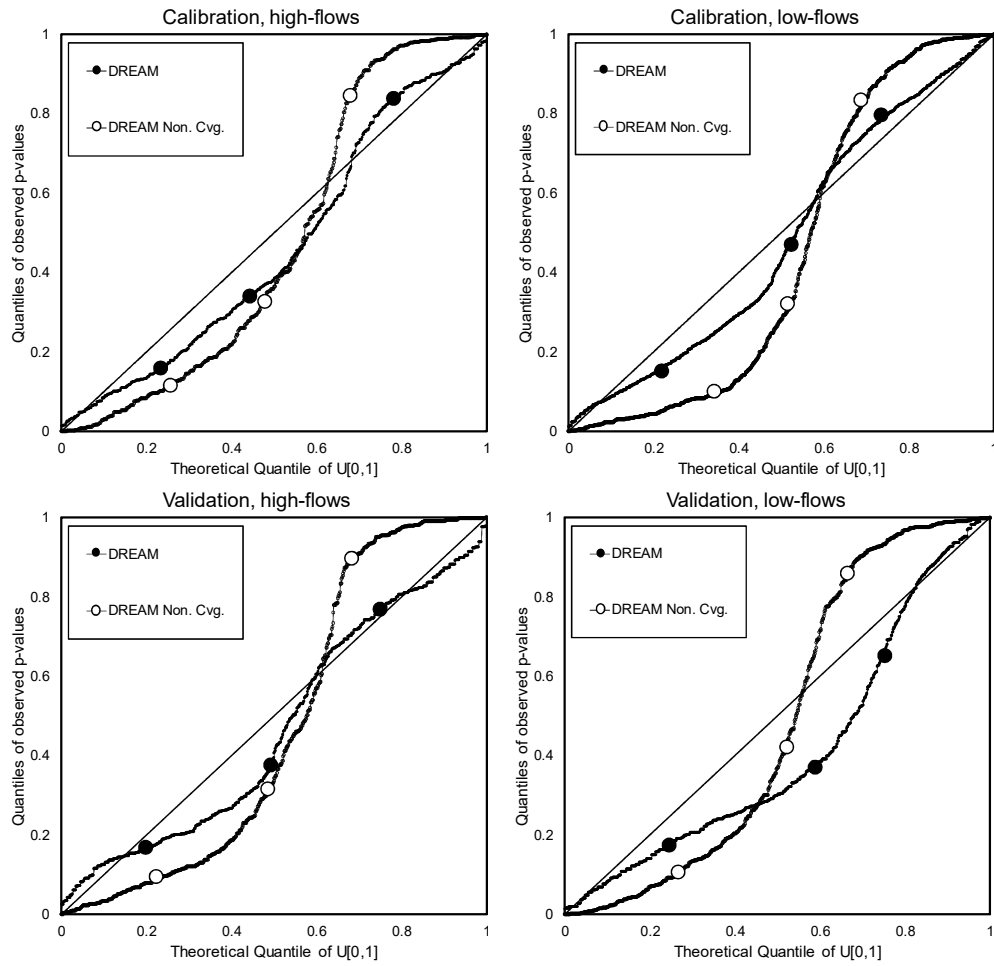
931 **Figure 7. Posterior ranges of WetSpa parameters derived by different uncertainty-based calibration techniques.**



932

933 **Figure 8. NS values of low-flows (horizontal axis) and high-flows (vertical axis) in calibration (upper panel) and validation (lower panel)**

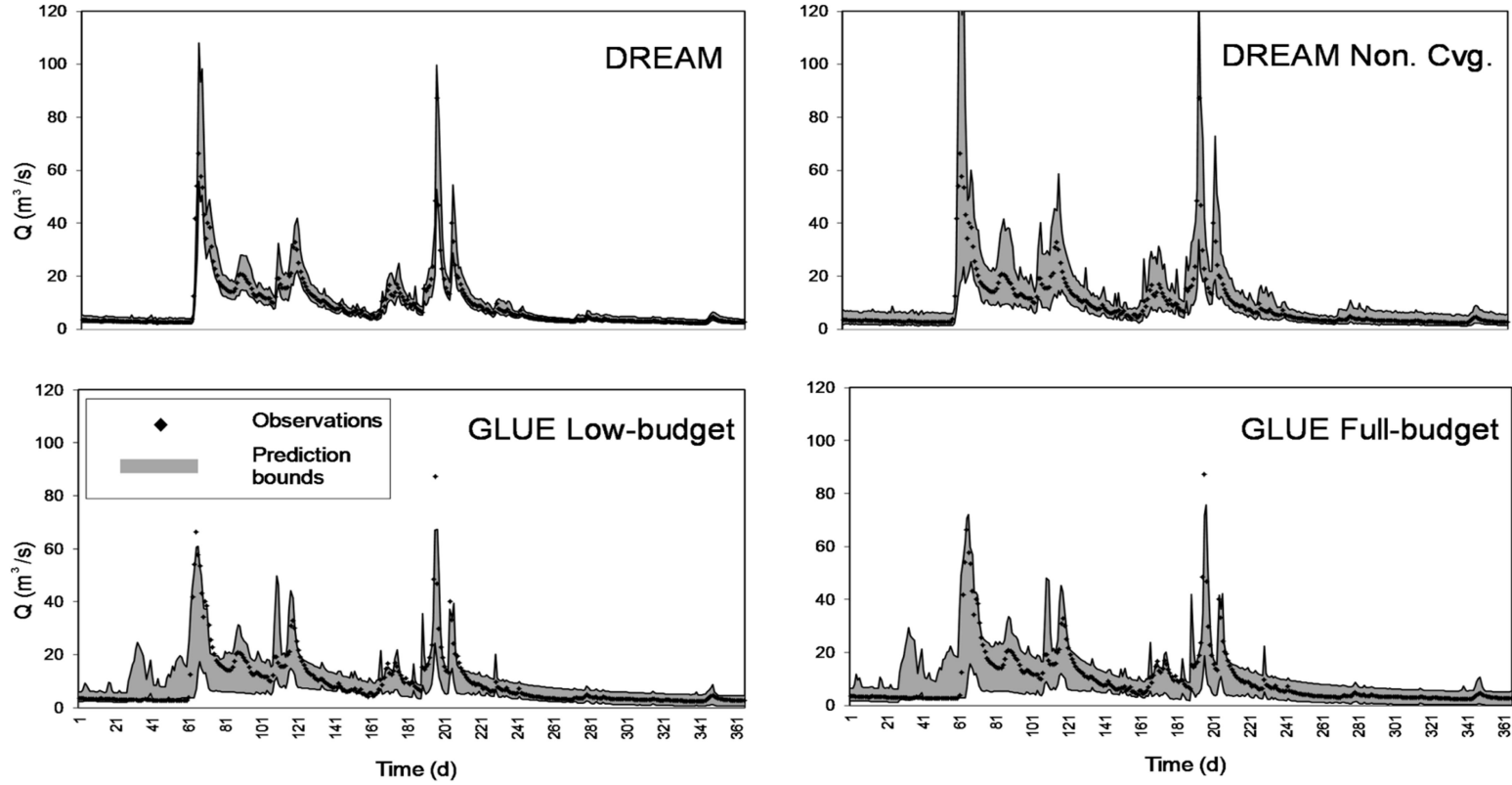
934 **period for WetSpa case study, derived from DREAM (light points) versus non-converged DREAM and GLUE methods (dark points).**



935

936 **Figure 9. QQ plot of Bayesian p-values for high- and low-flows derived from converged and non-**
 937 **converged DREAM, for the calibration (upper) and validation (bottom) periods of the WetSpa**
 938 **Hornad River case study.**

939



940

941 **Figure 10. Prediction bounds and observations for the year 1999 of validation period for the WetSpa Hornad River case study.**

942

943 List of Figures:

944 Figure 1. Schematic of the predictive QQ plot based on Thyer et al. (2009)

945 Figure 2. Posterior ranges of HYMOD parameters for the Leaf River case study; The parameter
946 ranges correspond to 95% posterior intervals for different uncertainty analysis methods.

947 Figure 3. NS values of low-flows (horizontal axis) and high-flows (vertical axis) in calibration
948 (upper panels) and validation (lower panels) period for HYMOD case study, derived from
949 DREAM (light points) versus non-converged DREAM and GLUE methods (dark points).

950 Figure 4. Validation period reliability and sharpness for low-flows (upper panels) and high-flows
951 (lower panels) in application of different techniques (shown in different shapes) to the HYMOD
952 and WetSpa simulation models.

953 Figure 5. QQ plot of Bayesian p-values for high- and low-flows derived from converged and
954 non-converged DREAM, for calibration (upper) and validation (bottom) periods of the HYMOD
955 Leaf River case study.

956 Figure 6. Prediction bounds and observations for the validation period in the HYMOD case
957 study.

958 Figure 7. Posterior ranges of WetSpa parameters derived by different uncertainty-based
959 calibration techniques.

960 Figure 8. NS values of low-flows (horizontal axis) and high-flows (vertical axis) in calibration
961 (upper panel) and validation (lower panel) period for WetSpa case study, derived from DREAM
962 (light points) versus non-converged DREAM and GLUE methods (dark points).

963 Figure 9. QQ plot of Bayesian p-values for high- and low-flows derived from converged and
964 non-converged DREAM, for the calibration (upper) and validation (bottom) periods of the
965 WetSpa Hornad River case study.

966 Figure 10. Prediction bounds and observations for the year 1999 of validation period for the
967 WetSpa Hornad River case study.