1

**Technical Note: Benefits of Bayesian estimation of model parameters in a large**

**hydrological model ensemble**

4

Yohei Sawada[1], and Shinichi Okugawa[1]

6

[1] Department of Civil Engineering, Graduate School of Engineering, the University of

Tokyo, Tokyo, Japan

9

10

Corresponding author: Y. Sawada, Department of Civil Engineering, the University of

Tokyo, Tokyo, Japan, 7-3-1, Hongo, Bunkyo-ku, Tokyo, Japan, yoheisawada@g.ecc.u-

tokyo.ac.jp

14

**Abstract**

Quantifying and mitigating parametric and structural uncertainties in hydrological models

are crucial to accurately understand and predict the rainfall-runoff process. Despite recent

advances in Bayesian approaches for quantifying structural uncertainty using very large

hydrological model ensembles, the simultaneous quantification of both parametric and

structural uncertainties has yet to be implemented since previous works on large model

ensembles have relied on deterministic optimization of parameters. Here we present the

potential benefits of Bayesian estimation of parametric uncertainty within a large

hydrological model ensemble. We find that Bayesian estimation of model parameters

(more generally, change in calibration methods) potentially influences the interpretation

25    of model comparisons. Specifically, Bayesian parametric uncertainty quantification

26    greatly benefits complex models with many parameters, thereby affecting discussions of

27    the appropriate level of model complexity. We also find that Bayesian parametric

28    uncertainty quantification does not substantially improve multi-model hydrological

29    predictions. The adverse effects of parameter misspecification in individual models are

30    effectively mitigated by combining models with diverse structures. Thus, the high

31    computational cost of Bayesian parameter estimation is not paid for to improve rainfall-

32    runoff analysis in a large hydrological model ensemble.

33

34

35

## 1. Introduction

36 **1. Introduction**

37 Hydrological models are essential tools to simulate the relationship between

38 meteorological conditions and runoff in river basins. Quantifying and mitigating

39 uncertainties in these models can greatly contribute to the accurate flood and drought

40 prediction, water resources management, and climate change assessment. Assuming

41 minimal error in input data, uncertainties in hydrological models can be broadly classified

42 into two categories: parametric uncertainty and structural uncertainty (e.g., Beven 2005;

43 Gupta et al. 2012). While parametric uncertainty arises from the misspecification of

44 model parameters, which are coefficients of equations represented in the model, structural

45 uncertainty originates from the specification of the equations themselves. The

46 quantification and mitigation of both parametric and structural uncertainties based on

47 hydrological observations, particularly river discharge data, remain grand challenges in

48 hydrology.

49

50 Mitigating parametric uncertainty has been extensively investigated in hydrology. Early

51 research largely focused on estimating a single set of parameters that minimized the cost

52 function measuring the fit between simulation and observation. In this paper, we refer

53 these approaches as deterministic optimization. A wide variety of gradient-based and

54 evolutionary algorithms have been applied to this task (e.g., Duan et al. 1993; Tolson et

55 al. 2007; Fowler et al. 2014; Qin et al. 2017, among others). Deterministic optimization

56 methods suffer from equifinality (Beven 2006), in which multiple combinations of

57 parameters reproduce observations equally well, and degrade the accuracy of simulating

58 unseen data. Bayesian estimation, by contrast, explicitly produces posterior probability

59 distributions based on observations and prior knowledge of parameters. This approach

60    offers several advantages: it can address equifinality (Vrugt et al. 2008a), enables

61    probabilistic forecasting (e.g., Hung et al. 2025), and consider errors in data such as

62    rainfall and river discharge measurements (e.g., Vrugt et al. 2008b). The golden standard

63    of Bayesian estimation is Markov Chain Monte Carlo (MCMC; Hastings 1970; Geman

64    and Geman 1984). In hydrology, the DiffeRential Evolution Adaptive Metropolis

65    (DREAM) algorithm (Vrugt et al. 2008c; Laloy and Vrugt 2012) has been widely used as

66    a variant of MCMC algorithms. While applying MCMC to hydrological models is

67    computationally expensive, sequential data assimilation offers a computationally cheap

68    alternative to estimate posterior distributions of model parameters and state variables (e.g.,

69    Moradkhani et al. 2005; Vrugt et al. 2013; Sawada 2022).

70

71    A common approach to mitigating structural uncertainty in process-based hydrological

72    models is the applications of machine-learning, which can provide fully data-driven

73    modeling (e.g., Kratzert et al. 2018, 2019) and correct the bias of process-based models

74    (e.g., Funato and Sawada 2025). These methods are intrinsically deterministic and remain

75    subject to equifinality. Large hydrological model ensembles, such as Modular Assessment

76    of Rainfall-Runoff Models Toolbox (MARRMoT; Knoben et al. 2019) and Raven (Craig

77    et al. 2020) opened the door to easily perform Bayesian uncertainty quantification of

78    model structure. Previous works employed Generalized Likelihood Uncertainty

79    Estimation (GLUE)-like methods (Beven and Binley 1992), in which "reasonably good"

80    models are selected based on predefined performance thresholds to identify equally

81    plausible models within a large multi-model ensemble. For instance, Knoben et al. (2020)

82    evaluated 36 hydrological models across 559 river basins using MARRMoT and

83    identified a high degree of structural equifinality. Knoben et al. (2025) further

84　demonstrated that this high equifinality can be mitigated by considering sampling

85　uncertainty in evaluation data. Chlumsky et al. (2021) performed simultaneous calibration

86　of model structures and parameters (see also Mai et al. 2020). They revealed a high degree

87　of equifinality within hydrological models implemented in the Raven framework.

88　Although these works on large samples of hydrological models have advanced the

89　quantification of structural uncertainty, they relied on deterministic parameter

90　optimization. Therefore, they did not explicitly consider parametric uncertainty in the

91　Bayesian way. It has yet to be clarified how Bayesian quantification of parametric

92　uncertainty benefits large hydrological model ensemble-based assessment of structural

93　uncertainty.

94

95　To address this research gap, we performed MCMC-based parametric uncertainty

96　quantification for 17 hydrological models with different structures across 51 river basins,

97　thereby enabling Bayesian estimation of both parametric and structural uncertainty. We

98　then examined scientific and practical benefits of applying Bayesian estimation of model

99　parameters in a large model ensemble. Specifically, we posed the following research

100　questions:

101　● Scientific benefits: Does Bayesian estimation of model parameters (or, more broadly,

102　　changes in calibration methods) potentially affect the interpretation of model

103　　comparisons?

104　● Practical benefits: Does Bayesian estimation of model parameters improve the

105　　accuracy of rainfall-runoff simulations in a large hydrological model ensemble? Is

106　　this improvement sufficient to justify the large computational cost of MCMC?

107

108

## 2. Method

### 2.1. Hydrological model

We used MARRMoT v1.3 (Knoben, 2019). MARRMoT includes 46 lumped conceptual

hydrological models with a wide range of complexities. The models are driven by daily

total precipitation, daily mean temperature, and daily mean potential evapotranspiration,

and they estimate daily basin-averaged runoff. From the 46 available models, we selected

models with IDs 02, 03, 04, 06, 07, 10, 11, 12, 13, 17, 18, 21, 24, 27, 30, and 31 (see the

supplement material of Knoben et al. 2019 for model details). These 17 models were

chosen to preserve diversity in model complexity while minimizing the total

computational cost, since Bayesian estimation of model parameters is computationally

expensive.

120

### 2.2. Parameter estimation

#### 2.2.1. Deterministic optimization

As a deterministic optimization method, we used the Nelder-Mead algorithm (Nelder and

Mead, 1965). We performed the deterministic parameter optimization for each model in

each basin. We used the *fminsearchbnd* function in MATLAB. Parameter ranges were

specified according to the original setting of MARRMoT (see Knoben et al. 2019).

Although our optimization method is a classical method and may be less capable of

avoiding local minima than modern evolutionary algorithms (e.g., Duan et al. 1993;

Tolson et al. 2007), we have found that it achieved high performance in our testbed

(Sawada et al. 2022; see Section 3).

131

### 2.2.2. Bayesian optimization

For Bayesian optimization, we adopted the method proposed by Liu et al. (2022). As the MCMC method, we used the DiffeRential Evolution Adaptive Metropolis (DREAM) algorithm (Vrugt et al. 2008c; Laloy and Vrugt 2012). While mean squared error is usually used as the formal likelihood function in the DREAM algorithm and many other MCMC applications, Liu et al. (2022) proposed using Kling-Gupta Efficiency (KGE: Gupta et al. 2009) as an informal likelihood function. Because KGE ranges from -infinity to one, Liu et al. (2022) applied a gamma density function to handle negative KGE values and developed a proper informal likelihood for the DREAM algorithm. We used the MATLAB implementation of DREAM ([https://github.com/Zaijab/DREAM](https://github.com/Zaijab/DREAM)). Besides the use of the KGE-based informal likelihood function, we used the default hyperparameter setting in this implementation. From the resulting Markov chains, we sampled 200 parameter sets to represent the posterior distribution of model parameters.

### 2.3. Bayesian model averaging

To evaluate the practical benefits of Bayesian parametric uncertainty quantification in a large hydrological model ensemble, we combined models with different structures and parameters into a single prediction using Bayesian Model Averaging (BMA).

First, we combined 17 models calibrated by deterministic optimization. Let the quantity of interest generated by the k-th model with deterministically optimized parameters, $M_k(\boldsymbol{\theta}_{DO,k})$, be denoted as $\Delta_k$, where $\boldsymbol{\theta}_{DO,k}$ represents the optimized parameters of the k-th model. In this study, $\Delta_k$ is runoff in the validation period (see Section 3). Given observation $\boldsymbol{y}$, the posterior mean of the quantity of interest is:

$$E[\Delta|\boldsymbol{y}] = \sum_{k}^{N} w_k \, \Delta_k \tag{1}$$

with model weights defined as:

$$w_k = P\big(M_k(\boldsymbol{\theta}_{DO,k})\big|y\big) = \frac{P\left(\boldsymbol{y}|M_k(\boldsymbol{\theta}_{DO,k})\right) P(M_k(\boldsymbol{\theta}_{DO,k})}{\sum_{i=1}^{N} P\left(\boldsymbol{y}|M_i(\boldsymbol{\theta}_{DO,i})\right) P(M_i(\boldsymbol{\theta}_{DO,i})} \tag{2}$$

where N is the total number of models (=17 in this case). $\boldsymbol{y}$ should be recognized as river discharge observation in a calibration period. Equation (1) indicates that the posterior mean is a weighted average of all model outputs, with weights proportional to their posterior probabilities. We parameterized the weights $w_k$ using KGE:

$$w_k = \frac{\exp\left(-\left(KGE_{max} - KGE\left(M_k(\boldsymbol{\theta}_{DO,k})\right)\right)\right)}{\sum_{i=1}^{N} \exp\left(-\left(KGE_{max} - KGE\left(M_i(\boldsymbol{\theta}_{DO,i})\right)\right)\right)} \tag{3}$$

where $KGE_{max}$ is the maximum KGE among the $N$ models, and $KGE(M_k(\boldsymbol{\theta}_{DO,k})$ is the KGE of the k-th model. In this paper, this averaging of hydrological models calibrated by deterministic optimization is specifically referred to as Bayesian Model Averaging (BMA).

Next, we compared the individual models calibrated by deterministic optimization with those calibrated by MCMC. To do so, we sampled 200 parameter sets from the posterior distributions, ran models with those parameter sets, and then combined these 200 hydrological simulations by Bayesian model averaging. Define the quantity of interest generated by the k-th model with the $l$-th parameter set, $M_k(\boldsymbol{\theta}_l)$, as $\Delta_{k,l}$. Similar to model averaging, the posterior mean of the quantity of interest estimated by the k-th model, $\Delta_k$ follows:

8

176
$$E[\Delta_k|\boldsymbol{y}] = \sum_{l}^{M} w_l \Delta_{k,l} \tag{4}$$

177 with weights:

178
$$w_l = \frac{\exp\left(-\left(KGE_{max,k} - KGE\left(M_k(\boldsymbol{\theta}_l)\right)\right)\right)}{\sum_{i=1}^{M} \exp\left(-\left(KGE_{max,k} - KGE\left(M_k(\boldsymbol{\theta}_i)\right)\right)\right)} \tag{5}$$

179 where $M$ is the total number of parameter sets (i.e. 200), $KGE_{max,k}$ is the maximum

180 KGE among the M parameter sets of the k-th model. This averaging within a single model

181 with different parameters is referred to as Bayesian Parameter Averaging (BPA).

182

183 Finally, we averaged all models and parameter sets. In this case, the posterior mean of the

184 quantity of interest is:

185
$$E[\Delta|\boldsymbol{y}] = \sum_{k}^{N} \sum_{l}^{M} w_{k,l} \Delta_{k,l} \tag{6}$$

186 with weights:

187
$$w_{k,l} = \frac{\exp\left(-\left(KGE_{max,k,l} - KGE\left(M_k(\boldsymbol{\theta}_l)\right)\right)\right)}{\sum_{i=1}^{N} \sum_{j=1}^{M} \exp\left(-\left(KGE_{max,k,l} - KGE\left(M_i(\boldsymbol{\theta}_j)\right)\right)\right)} \tag{7}$$

188 where $KGE_{max,k,l}$ is the maximum KGE among all models and parameter combinations.

189 This joint averaging is called Bayesian Model and Parameter Averaging (BMPA).

190

191 **3. Experiment design**

192 We applied the aforementioned methods to 51 river basins in Japan (Figure 1). We used

193 meteorological forcings from the Multi-model Ensemble for Robust Verification of

194 hydrological modeling in Japan (MERV-Jp) dataset (Sawada et al. 2022; Sawada and

195  Okugawa 2022). The river basins shown in Figure 1 cover a wide range of climatic, soil,

196  land-use, anthropogenic, and topographic conditions. Sawada et al. (2022) reported that

197  44 deterministically calibrated models in MARRMoT achieved high accuracy to

198  reproduce observed runoff. The best KGEs in 44 models exceeded 0.7 in nearly all basins,

199  which are comparable to those reported in other large model ensemble studies across

200  different regions (e.g., Knoben et al. (2020)). Therefore, our findings are expected to be

201  transferable to similar studies in the context of large hydrological model ensembles.

202

203  The study period spans 1986-2015. Two calibration periods were considered: a 5-year

204  calibration period representing data-rich conditions, and a 1-year calibration period

205  representing data-poor conditions. In the 5-year calibration scenario, the initial 5-year

206  (1986-1990) data were used for calibration with both deterministic optimization and

207  MCMC, and the remaining 25-year (1991-2015) data were used for evaluation. In basins

208  where complete discharge records were unavailable for 1986-1990, the calibration period

209  was shifted to ensure a continuous 5-year record, which slightly reduced the validation

210  period. We used the same 5-year data for model spin-up, resulting a 10-year model

211  integration for each parameter evaluation step. In the 1-year calibration scenario, we

212  applied the same 5-year data spin-up, followed by evaluation of parameters in the

213  subsequent 1-year simulation (1986 in most basins). The validation data for the 1-year

214  calibration scenario were identical to those for the 5-year calibration scenario. Model

215  performance was evaluated using KGE and Nash-Sutcliffe Efficiency (NSE; Nash and

216  Sutcliffe 1970; see also Duc and Sawada 2023 for a modern interpretation of NSE) during

217  the validation period.

218

219  For deterministic optimization, we have 17 hydrological predictions by 17 models in 51

220  river basins, yielding 17×51 = 867 hydrographs. For Bayesian estimation of parametric

221  uncertainty, each of the 17 models was run with 200 posterior parameter sets, producing

222  17×200 = 3400 simulations per basin. Across 51 basins, this amounted to 3400×51 =

223  173,400 hydrographs. First, we compared the performance of individual models

224  calibrated by deterministic optimization with those calibrated by MCMC and combined

225  by BPA (see Section 2.3) across all 17 models in 51 river basins, to assess if the evaluation

226  of the model structures is affected by parameter estimation method. Second, we compared

227  the performance of BMA and BPMA (see Section 2.3), to discuss the potential benefits

228  of considering parametric uncertainty through Bayesian estimation in improving

229  hydrological predictions.

230

231  **4. Results and discussions**

232  Bayesian estimation of model parameters (i.e., BPA) systematically outperforms

233  deterministic optimization. Boxplots of Figure 2 show the differences in KGE and NSE

234  between BPA and deterministically optimized models. The superiority of MCMC-based

235  optimization is consistent with earlier findings on the DREAM (e.g., Laloy and Vrugt

236  2012) algorithm. To our best knowledge, we verified this superiority within a large model

237  ensemble for the first time. Two main factors explain this result. First, MCMC explores

238  the parameter space more extensively to maximize KGE than deterministic optimization.

239  Second, MCMC accounts for equifinality by sampling multiple parameter sets that

240  reproduce observation equally well, which leads the higher robustness to unseen data.

241  Even in the rich-data scenario (i.e., the 5-year calibration), improvements in KGE exceeds

242  0.2 in some models. This trend is more pronounced in the data-scarce scenario (i.e., 1-

243    year calibration scenario). When calibration data are limited, parameter estimates are

244    inherently uncertain, since many different parameter combinations may be able to equally

245    explain the limited data. In such cases, Bayesian estimation is more appropriate than

246    deterministic methods.

247

248    Figure 3 reveals that the improvements achieved by BPA over deterministically optimized

249    models systematically appear. For instance, model ID 10 (Susannah Brook model v2; see

250    Son and Sivapalan 2007, Knoben et al. 2019) shows substantial gains from Bayesian

251    parameter estimation in more than 20 river basins. Except for this model, more complex

252    models with the larger numbers of parameters (e.g., ID 30, 31, and 32) tend to benefit

253    more from Bayesian estimation than simpler models (e.g., ID 2, 3, and 4) (note that

254    models with higher IDs generally correspond to greater structural complexity; see

255    Knoben et al. 2019) especially in the data-rich scenario (i.e., 5-year calibration scenario).

256    This indicates that models with many parameters are particularly affected by equifinality

257    and therefore gain substantially from Bayesian parameter estimation.

258

259    Previous studies have evaluated model structures by comparing the performance of

260    deterministically optimized models. Our results imply that such evaluations may be

261    affected by the choice of parameter optimization methods. For instance, although results

262    in Knoben et al. (2020) and Knoben et al. (2025) indicated that complex models with

263    many parameters do not necessarily outperform simpler models, this conclusion may

264    partly reflect an underestimation of the maximum potential performance of complex

265    models due to reliance on deterministic optimization. While these complex models have

266    been shown not to suffer from overfitting (Knoben et al. 2020), there might be room for

267  improvement through Bayesian estimation, which explicitly addresses parameter

268  equifinality.

269

270  When a large number of calibrated models are available, a practical way to improve

271  prediction accuracy is to use the (weighted) average of their outputs (e.g., Kimizuka and

272  Sawada 2022; Zhang and Yang 2018). The red dots in Figure 2 show the performance

273  differences between BPMA and BMA. Although Bayesian parameter estimation

274  substantially improves the performance of individual models, the overall improvement

275  from BMA to BPMA is marginal. Even in the data-scarce scenario (i.e., the 1-year

276  calibration scenario), the improvement in KGE (NSE) by Bayesian estimation is less than

277  0.1 (0.05). Surprisingly, even the classic optimization method remains competitive with

278  DREAM-based Bayesian optimization when models are combined in a large ensemble.

279  Considering the substantial computational costs of the MCMC-based Bayesian parameter

280  estimation, we conclude that Bayesian parametric uncertainty quantification provides

281  limited practical benefits for improving predictions in large hydrological model

282  ensembles.

283

284  This occurs because poorly performing models produced by deterministic optimization

285  are assigned lower weights in the BMA framework. Figure 4 illustrates a typical case: in

286  basin no. 43, BPA achieves 0.8 KGE for nearly all models, while some deterministically

287  optimized models (i.e., ID 7, 24, 30, and 31) perform poorly. Nevertheless, BMA remains

288  competitive with BPMA, since the poorly performing models receive smaller weights

289  during averaging. Therefore, the adverse effects of suboptimal calibration are effectively

290  mitigated.

291

**5. Conclusions**

293 Here we performed MCMC-based parameter optimization for 17 hydrological models

294 across 51 river basins to clarify the potential benefits of Bayesian parametric uncertainty

295 quantifications in a large hydrological model ensemble. Scientifically, Bayesian

296 parametric uncertainty quantification is important because it can influence the

297 interpretation of structural uncertainty assessment. The benefits of the Bayesian

298 parameter estimation appear systematically across models rather than randomly. Certain

299 models gain substantial improvements, and more complex models tend to benefit more

300 strongly than simpler models. Considering Bayesian parameter uncertainty potentially

301 affects the discussion of the appropriate complexity of hydrological models.

302

303 Practically, Bayesian parametric uncertainty quantification does not greatly contribute to

304 improving multi-model ensemble hydrological prediction. It implies that structural errors

305 are larger than parametric errors in hydrological prediction. Given the high computational

306 cost of Bayesian estimation, multi-model ensembles calibrated by deterministic

307 optimization are sufficient in many cases.

308

309 Our analysis was limited to 17 models, fewer than in previous large ensemble studies

310 (Knoben et al. 2020, 2025; Chlumsky et al. (2021)). We had to exclude computationally

311 expensive models in MARRMoT to make the MCMC applications feasible in this initial

312 attempt. Future work should expand to all MARRMoT models and pursue GLUE-like

313 assessments of structural uncertainty using Bayesian parameter uncertainty quantification

314 by unleashing the power of high-performance computers.

315

**Acknowledgement**

322

**Code availability**

MARRMoT v1.3 is available at https://doi.org/10.5281/zenodo.3235664 (Knoben, 2019).

DREAM is available at https://github.com/Zaijab/DREAM.

326

**Data availability**

Results of hydrological models in this work can be found at https://doi.org/10.5281/zenodo.17282833 (Sawada and Okugawa 2025).

330

**Author contribution**

YS designed the study, interpreted results, and wrote the initial version of the paper. SO performed numerical experiments, analyzed the results, and contributed to editing the paper.

335

**References**

Beven, K.: On the concept of model structural error, Water Sci. Technol., 52, 167–175, https://doi.org/10.2166/wst.2005.0165, 2005.

339    Beven, K.: A manifesto for the equifinality thesis, J. Hydrol., 320, 18–36,

340    https://doi.org/10.1016/j.jhydrol.2005.07.007, 2006.

341    Chlumsky, R., Mai, J., Craig, J. R., and Tolson, B. A.: Simultaneous calibration of

342    hydrologic model structure and parameters using a blended model, Water Resour. Res.,

343    57, https://doi.org/10.1029/2020WR029229, 2021.

344    Craig, J. R., Brown, G., Chlumsky, R., Jenkinson, R. W., Jost, G., Lee, K., Mai, J., Serrer,

345    M., Sgro, N., Shafii, M., Snowdon, A. P., and Tolson, B. A.: Flexible watershed simulation

346    with the Raven hydrological modelling framework, Environ. Model. Softw., 129, 104728,

347    https://doi.org/10.1016/j.envsoft.2020.104728, 2020.

348    Duan, Q. Y., Gupta, V. K., and Sorooshian, S.: Shuffled complex evolution approach for

349    effective and efficient global minimization, J. Optim. Theory Appl., 76, 501–521,

350    https://doi.org/10.1007/BF00939380, 1993.

351    Duc, L. and Sawada, Y.: A signal-processing-based interpretation of the Nash–Sutcliffe

352    efficiency, Hydrol. Earth Syst. Sci., 27, 1827–1839, https://doi.org/10.5194/hess-27-

353    1827-2023, 2023.

354    Fowler, K., Peel, M., Western, A., and Zhang, L.: Improved rainfall-runoff calibration for

355    drying climate: Choice of objective function, Water Resour. Res., 54, 3392–3408,

356    https://doi.org/10.1029/2017WR022466, 2018.

357    Funato, M. and Sawada, Y.: Multi-model ensemble and reservoir computing for river

358    discharge     prediction      in      ungauged      basins,      arXiv      [preprint],

359    https://arxiv.org/abs/2507.18423, 2025.

360    Geman, S. and Geman, D.: Stochastic relaxation, Gibbs distributions and the Bayesian

361    restoration of images, IEEE Trans. Pattern Anal. Mach. Intell., 6, 721–741, 1984.

362   Gupta, H. V., Clark, M. P., Vrugt, J. A., Abramowitz, G., and Ye, M.: Towards a

363   comprehensive assessment of model structural adequacy, Water Resour. Res., 48,

364   https://doi.org/10.1029/2011WR011044, 2012.

365   Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean

366   squared error and NSE performance criteria: Implications for improving hydrological

367   modelling, J. Hydrol., 377, 80–91, https://doi.org/10.1016/j.jhydrol.2009.08.003, 2009.

368   Hastings, W. K.: Monte Carlo sampling methods using Markov chains and their

369   applications, Biometrika, 57, 97–109, https://doi.org/10.1093/biomet/57.1.97, 1970.

370   Kimizuka, T. and Sawada, Y.: How robust is a multi-model ensemble mean of conceptual

371   hydrological   models   to   climate   change?,   Water,   14,   2852,

372   https://doi.org/10.3390/w14182852, 2022.

373   Knoben,   W.   J.   M.:   MARRMoT_v1.2,   Zenodo   [code],

374   https://doi.org/10.5281/zenodo.3235664, 2019.

375   Knoben, W. J. M., Freer, J. E., Fowler, K. J. A., Peel, M. C., and Woods, R. A.: Modular

376   Assessment of Rainfall–Runoff Models Toolbox (MARRMoT) v1.2: an open-source,

377   extendable framework providing implementations of 46 conceptual hydrologic models as

378   continuous   state-space   formulations,   Geosci.   Model   Dev.,   12,   2463–2480,

379   https://doi.org/10.5194/gmd-12-2463-2019, 2019.

380   Knoben, W. J. M., Freer, J. E., Peel, M. C., Fowler, K. J. A., and Woods, R. A.: A brief

381   analysis of conceptual model structure uncertainty using 36 models and 559 catchments,

382   Water Resour. Res., 56, e2019WR025975, https://doi.org/10.1029/2019WR025975, 2020.

383   Knoben, W. J. M., Raman, A., Gründemann, G. J., Kumar, M., Pietroniro, A., Shen, C.,

384   Song, Y., Thébault, C., van Werkhoven, K., Wood, A. W., and Clark, M. P.: Technical

385   note: How many models do we need to simulate hydrologic processes across large

386   geographical    domains?,    Hydrol.    Earth    Syst.    Sci.,    29,    2361–2375,

387   https://doi.org/10.5194/hess-29-2361-2025, 2025.

388   Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall–runoff

389   modelling using Long Short-Term Memory (LSTM) networks, Hydrol. Earth Syst. Sci.,

390   22, 6005–6022, https://doi.org/10.5194/hess-22-6005-2018, 2018.

391   Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., and Nearing, G.:

392   Towards learning universal, regional, and local hydrological behaviors via machine

393   learning applied to large-sample datasets, Hydrol. Earth Syst. Sci., 23, 5089–5110,

394   https://doi.org/10.5194/hess-23-5089-2019, 2019b.

395   Laloy, E. and Vrugt, J. A.: High-dimensional posterior exploration of hydrologic models

396   using multiple-try DREAM(ZS) and high-performance computing, Water Resour. Res.,

397   48, W01526, https://doi.org/10.1029/2011WR010608, 2012.

398   Liu, Y., Fernández-Ortega, J., Mudarra, M., and Hartmann, A.: Pitfalls and a feasible

399   solution for using KGE as an informal likelihood function in MCMC methods:

400   DREAM(ZS)    as    an    example,    Hydrol.    Earth    Syst.    Sci.,    26,    5341–5355,

401   https://doi.org/10.5194/hess-26-5341-2022, 2022.

402   Mai, J., Craig, J. R., and Tolson, B. A.: Simultaneously determining global sensitivities

403   of model parameters and model structure, Hydrol. Earth Syst. Sci., 24, 5835–5858,

404   https://doi.org/10.5194/hess-24-5835-2020, 2020.

405   Moradkhani, H., Hsu, K.-L., Gupta, H., and Sorooshian, S.: Uncertainty assessment of

406   hydrologic model states and parameters: Sequential data assimilation using the particle

407   filter, Water Resour. Res., 41, https://doi.org/10.1029/2004WR003604, 2005.

408    Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models. Part 1:

409    A discussion of principles, J. Hydrol., 10, 282–290, https://doi.org/10.1016/0022-

410    1694(70)90255-6, 1970.

411    Nelder, J. A. and Mead, R.: A simplex method for function minimization, Comput. J., 7,

412    308–313, https://doi.org/10.1093/comjnl/7.4.308, 1965.

413    Qin, Y., Kuczera, G., and Kavetski, D.: Comparison of Newton-type and SCE

414    optimisation algorithms for the calibration of conceptual hydrological models, Australas.

415    J. Water Resour., 20, 169–176, https://doi.org/10.1080/13241583.2017.1298180, 2016.

416    Sawada, Y.: An efficient estimation of time-varying parameters of dynamic models by

417    combining offline batch optimization and online data assimilation, J. Adv. Model. Earth

418    Syst., 14, e2021MS002882, https://doi.org/10.1029/2021MS002882, 2022.

419    Sawada, Y., Okugawa, S., and Kimizuka, T.: Multi-model ensemble benchmark data for

420    hydrological modeling in Japanese river basins, Hydrol. Res. Lett., 16, 73–79,

421    https://doi.org/10.3178/hrl.16.73, 2022.

422    Sawada, Y., and Okugawa, S.: Dataset of "Benefits of Bayesian estimation of model

423    parameters    in    a    large    hydrological    model    ensemble"    [Dataset],

424    https://doi.org/10.5281/zenodo.17282833, 2025

425    Son, K. and Sivapalan, M.: Improving model structure and reducing parameter

426    uncertainty in conceptual water balance models through the use of auxiliary data, Water

427    Resour. Res., 43, W01415, https://doi.org/10.1029/2006WR005032, 2007.

428    Tolson, B. A. and Shoemaker, C. A.: Dynamically dimensioned search algorithm for

429    computationally efficient watershed model calibration, Water Resour. Res., 43,

430    https://doi.org/10.1029/2005WR004723, 2007.

431    Vrugt, J. A., Ter Braak, C. J. F., Clark, M. P., Hyman, J. M., and Robinson, B. A.:

432    Treatment of input uncertainty in hydrologic modeling: Doing hydrology backward with

433    Markov chain Monte Carlo simulation, Water Resour. Res., 44,

434    https://doi.org/10.1029/2007WR006720, 2008b.

435    Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Robinson, B. A., and Hyman, J. M.:

436    Accelerating Markov chain Monte Carlo simulation by differential evolution with self-

437    adaptive randomized subspace sampling, Int. J. Nonlinear Sci. Numer. Simul., 10, 273–

438    290, 2008c.

439    Vrugt, J. A., ter Braak, C. J. F., Gupta, H. V., and Robinson, B. A.: Equifinality of formal

440    (DREAM) and informal (GLUE) Bayesian approaches in hydrologic modeling?, Stoch.

441    Environ. Res. Risk Assess., 23, 1011–1026, https://doi.org/10.1007/s00477-008-0274-y,

442    2008a.

443    Vrugt, J. A., Ter Braak, C. J. F., Diks, C. G. H., and Schoups, G.: Hydrologic data

444    assimilation using particle Markov chain Monte Carlo simulation: Theory, concepts and

445    applications, Adv. Water Resour., 51, 457–478,

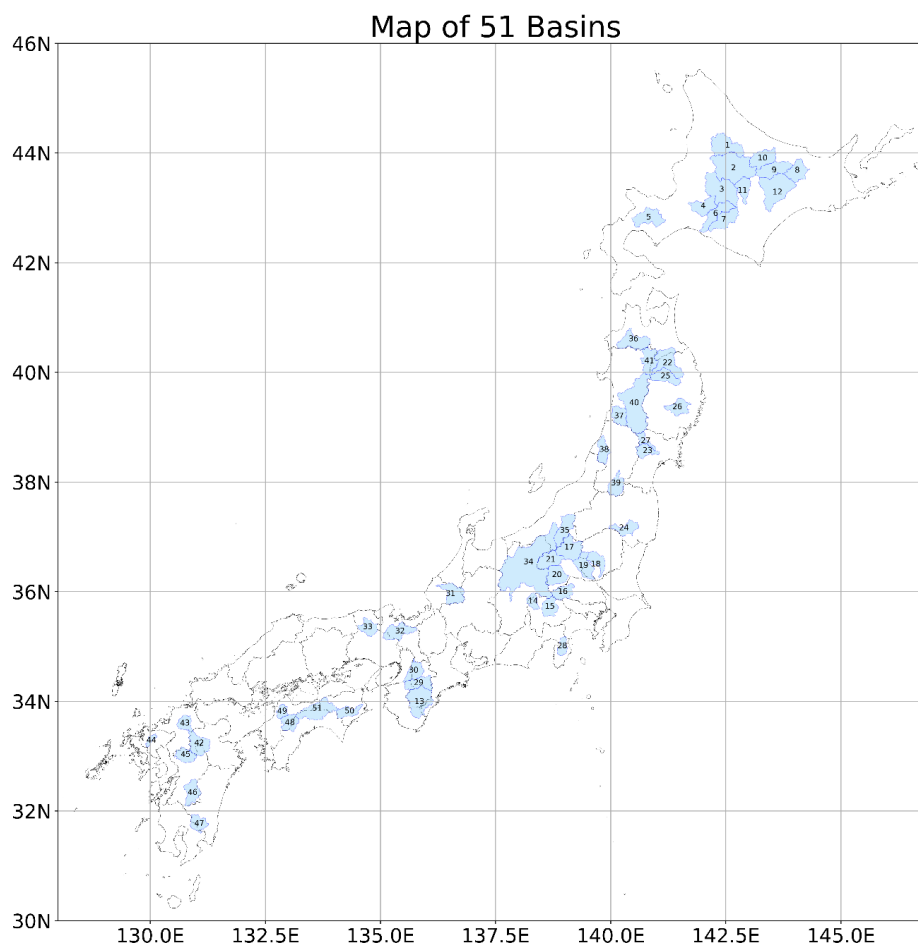446    https://doi.org/10.1016/j.advwatres.2012.04.002, 2013.

447    Zhang, L. and Yang, X.: Applying a multi-model ensemble method for long-term runoff

448    prediction under climate change scenarios for the Yellow River Basin, China, Water, 10,

449    301, https://doi.org/10.3390/w10030301, 2018.

450

451

452



453

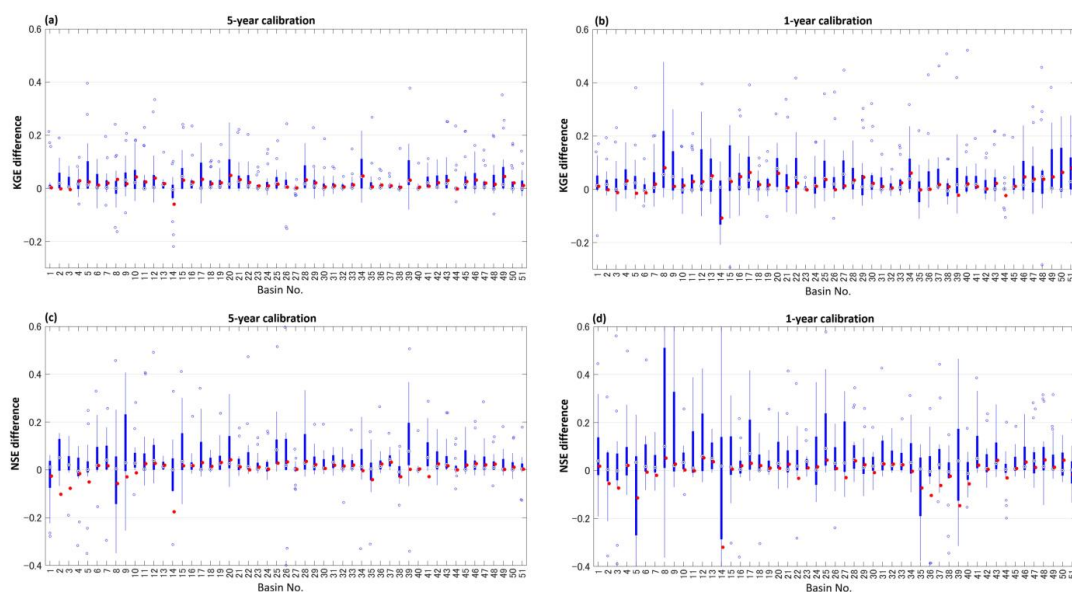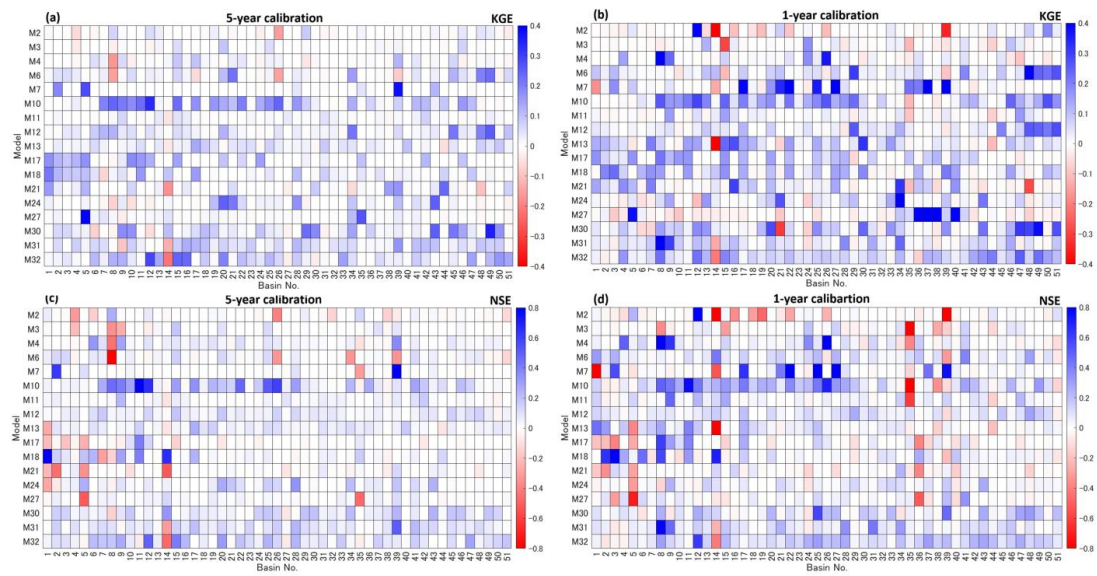**Figure 1**. Study area of 51 river basins.

455

456

**Figure 2**. Differences in KGE (a, b) and NSE (c, d) between BPA and deterministically optimized models (boxplots) in 5-year (a, c) and 1-year (b, d) calibration scenarios. Red dots show the performance differences between BMPA and BMA (see Section 3).
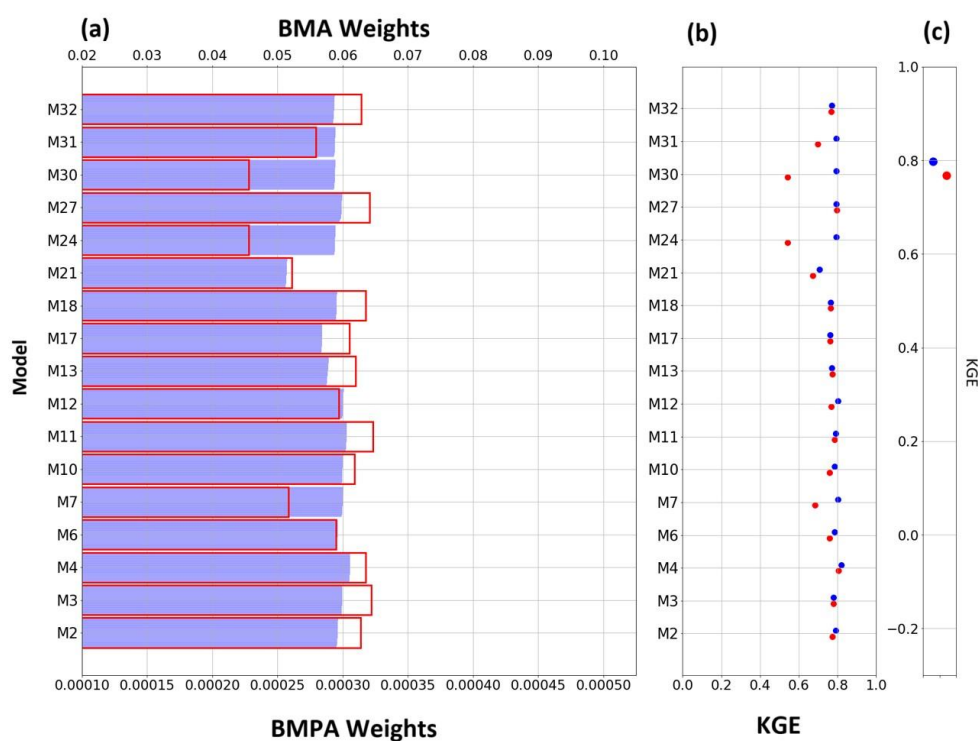
462



463

**Figure 3**. Differences in KGE (a, b) and NSE (c, d) between BPA and deterministically optimized models for

each basin and model in 5-year (a, c) and 1-year (b, d) calibration scenarios.

**Figure 4.** (a) Weights in BMA (red bars) and BMPA (blue bars). (b) KGE of BPA (blue dots) and deterministically optimized models (red dots). (c) KGE of BMPA (blue dots) and BMA (red dots). Results of basin No. 43 (see Figure 1) are shown.