

# Quantifying input uncertainty in the calibration of water quality models: reshufflingreordering errors via the secant method

Xia Wu<sup>1,2</sup>, Lucy Marshall<sup>2</sup>, Ashish Sharma<sup>2</sup>

<sup>1</sup>College of Hydrology and Water Resources, Hohai University, Nanjing, 210098, China

5 <sup>2</sup>School of Civil and Environmental Engineering, University of New South Wales, Sydney, 2052, Australia

*Correspondence to:* Lucy Marshall (lucy.marshall@unsw.edu.au)

**Abstract.** Uncertainty in inputs can significantly impair parameter estimation in water quality modeling, necessitating accurate quantification of input errors. However, decomposing input error from model residual error is still challenging. This study develops a new algorithm, referred to as Bayesian error analysis with reshufflingreordering (BEAR), to address this problem.

10 The basic approach requires sampling errors from a pre-estimated error distribution and then reshufflingreordering them with their inferred ranks via the secant method. This approach is demonstrated in the case of total suspended solids (TSS) simulation via a conceptual water quality model. Based on case studies using synthetic data, the BEAR method successfully improves the identification of the input errors in the model calibration. The results of a real case study demonstrate that even with the presence of model structural error and output data error, the BEAR method can approximate the true input and bring a better  
15 model fit through an effective input modification. However, its effectiveness is limited by the ~~assumption that the input uncertainty should be dominant~~accuracy and ~~that the prior information~~selection of the input error model ~~can be estimated.~~. The application of the BEAR method in TSS simulation ~~is effective for understanding a range of water quality conditions and the further developed algorithm~~ can be extended to other water quality predictions~~models~~.

## 1 Introduction

20 For robust water management, uncertainty analysis is of growing importance in water quality modeling (Refsgaard et al., 2007). It can provide knowledge of error propagation and the magnitude of uncertainty impacts in model simulations to guide improved predictive performance (Radwan et al., 2004). However, the implementation of uncertainty analysis in water quality models (WQMs) is still challenging due to complex interactions among sources of multiple errors, generally caused by a  
25 simplified model structure (structural uncertainty), imperfect observed data (input uncertainty and observation uncertainty in calibration data) and limited parameter identifiability (parametric uncertainty) (Refsgaard et al., 2007).

Among them, input uncertainty is expected to be particularly significant in a WQM, interpreted here as the observation uncertainty of any input data. Observation uncertainty is different from other sources of uncertainty in modeling since these uncertainties arise independently of the WQM itself, thus, their properties (e.g. probability distribution family and distribution parameters) can, at least in principle, be estimated prior to the model calibration and simulation by analysis of the data

30 acquisition instruments and procedures (McMillan et al., 2012). Rode and Suhr (2007) and Harmel et al. (2006) reviewed the uncertainty associated with selected water quality variables based on the empirical quality of observations. The general methodology developed in their studies can be extended to the analysis of other water quality variables. Besides the error coming from the measurement process, the error from surrogated data is another major source of input uncertainty (McMillan et al., 2012). Measurements of water quality variables often lack desirable temporal and spatial resolutions, thus, the use of  
35 surrogate or proxy data is necessary for improved inference of water quality parameters (Evans et al., 1997, Stubblefield et al., 2007). For the surrogate error, its probability distribution is easy to estimate from the residuals between the measurements and proxy values. In this process, the measurement errors are ignored given the errors introduced from the surrogate process are commonly much more than the measurement errors (McMillan et al., 2012). These estimated error distributions are “prior knowledge” of input uncertainty before any model calibration and can serve as the a-priori uncertainty estimation in the  
40 modeling process.

Input uncertainty can lead to bias in parameter estimation in water quality modeling (Chaudhary and Hantush, 2017, Kleidorfer et al., 2009, Willems, 2008). Improved model calibration requires isolating the input uncertainty from the total uncertainty. However, the precise quantification of time-varying input errors is still challenging when other types of uncertainties are propagated through to the model results. In hydrological modeling, several approaches have been developed to characterize  
45 time-varying input errors, and these may hold promise for application in WQMs. The Bayesian total error analysis (BATEA) method provides a framework that has been widely used (Kavetski et al., 2006). Time-varying input errors are defined as multipliers on the input time series and inferred along with the model parameters in a Bayesian calibration scheme. This leads to a high-dimensionality problem, which cannot be avoided (Renard et al., 2009) and restricts the application of this approach to the assumption of event-based multipliers (the same multiplier applied to one storm event). In the Integrated Bayesian  
50 Uncertainty Estimator (IBUNE) (Ajami et al., 2007) approach, multipliers are not jointly inferred with the model parameters, but sampled from the assumed distribution and then filtered by the constraints of simulation fitting. This approach reduces the dimensionality significantly and can be applied in the assumption of data-based multiplier (one multiplier for one input data) (Ajami et al., 2007). However, this approach results is less effective because the probability of co-occurrence of all optimal error/parameter values is very low, resulting in an underestimation of the multiplier variance and misidentification of the  
55 uncertainty sources (Renard et al., 2009). From the above, a new strategy should be developed to avoid high dimensional computation and ensure the accuracy of error identification.

To complete this goal, this study develops a new algorithm – Bayesian error analysis with reshufflingreordering (BEAR). The derivation and details of the BEAR algorithm in quantifying input errors are described in Sect. 2. Section 3 introduces the build-up/wash-off model (BwMod) to illustrate this approach. Its model input, streamflow, often suffers from observational  
60 errors from a rating curve. By comparing the results with other calibration frameworks, the ability of the BEAR method is explored in atwo synthetic easecases and a real case. In this way, the new algorithm is tested in a simplecontrolled situation (with an assumptionthe knowledge of the true outputerror and data and model structurevalue) and in a realistic situation (with

the interference of multiple error sources) respectively. Section 4 evaluates the BEAR method and its implementation. Finally, Section 5 outlines the main conclusions and recommendations for this work.

## 65 2 Methodology

### 2.1 Basic theory of identifying the input error in model calibration

A WQM in the ideal situation without any error can be described as

$$Y^* = M(X^* | \theta^*) \quad (1)$$

where the ~~true~~ asterisk \* implies the true value without error, and the true output  $Y^*$  is simulated by the perfect model  $M$  with the true input  $X^*$  and the true model parameter  $\theta^*$ . Here and in the following contents, a capital bold letter (e.g.  $X, Y$ ) represents a vector and a lower case (e.g.  $x, y$ ) represents a variablesscalar.

In reality, the model input  $X^o$  (typically the rainfall or streamflow in a WQM) inevitably suffers from input error  $\epsilon_x$ . This will result in a calibrated model parameter  $\theta^c$  biased from the true value  $\theta^*$  (Kleidorfer et al., 2009). Thus, under the assumption that the output data and model structure are generally without errors and the input errors are additive to the true input data  $X^*$ , the model residual  $\epsilon$  in a traditional calibration can be described by

$$\epsilon = Y^o - Y^s = Y^o - M(X^o | \theta^c) = Y^* - M(X^* + \epsilon_x | \theta^c) \quad (2)$$

Under the ideal situation without input errors, the residual will reduce to zero, like

~~$$\epsilon = Y^o - Y^s = Y^* - M(X^* | \theta^*) = 0 \quad (3)$$~~

where  $Y^s$  is the output simulated from the model  $M$  corresponding to the observed input  $X^o$  and model parameter  $\theta^c$ , and the observed output  $Y^o$  is assumed without observational errors in the derivation, thus can be denoted as  $Y^*$ .

It should be noted that the derivation of the BEAR method is based on the assumption that the model only suffers from input error and parameter error, but other sources of error (i.e. model structural error and output observational error) can also impair the estimation of the model parameters and are inevitable in the WQM. Considering this realistic situation, the ability of the BEAR method will be tested in a case study where the interference of other sources of error has been considered.

To counter the influence of input errors in a traditional calibration, an appealing approach is to subtract estimated errors  $\epsilon_x^p$  from the observed input  $X^o$ . This is illustrated as the “proposed” approach and the superscript  $p$  represents the values in this “proposed” approach. The residual  $\epsilon^p$  will change to

$$\boldsymbol{\varepsilon}^p = \mathbf{Y}^o - \mathbf{Y}^p = \mathbf{Y}^* - M(\mathbf{X}^p | \boldsymbol{\theta}^p) = \mathbf{Y}^* - M(\mathbf{X}^* + \boldsymbol{\varepsilon}_x - \boldsymbol{\varepsilon}_x^p | \boldsymbol{\theta}^p) \quad (3)$$

If the equivalence between  $\boldsymbol{\varepsilon}_x$  and  $\boldsymbol{\varepsilon}_x^p$  can be ensured for each data point, the modified input  $\mathbf{X}^p$  then becomes the same as the true value  $\mathbf{X}^*$ . The proposed calibration (Eq. (3)) will turn into an ideal calibration where the optimal parameters  $\boldsymbol{\theta}^p$  will lead to the same simulation corresponding to the true values  $\boldsymbol{\theta}^*$  and the model residual  $\boldsymbol{\varepsilon}^p$  will decrease to zero. If the inverse problem (from the zero residual to find the optimal parameter) is not unique, the calibrated parameter  $\boldsymbol{\theta}^p$  may not converge to the true parameter  $\boldsymbol{\theta}^*$ , but lead to the same simulation as the true parameter. In this study, these parameters are also denoted as  $\boldsymbol{\theta}^*$  and called ideal model parameters. Besides, if the identified input error and the model parameter can compensate each other, multiple combinations of model parameter and input error may yield zero residual and their estimates will be biased from the ideal values. A possible way to weaken this compensation effect will be explored Sect. 4.2. Although the aforementioned problems cannot be avoided, selecting the optimal input error series according to the model residual error is the basic theory of not only this study but also current methods identifying the input errors (i.e. BATEA (Kavetski et al., 2006) and IBUNE (Ajami et al., 2007)).

The above approach does not improve the input error model itself but improves the WQM specification to have parameters closer to what would be achieved under no error conditions. Then the model can be more effectively used for scenario analysis (where we may know the hydrologic regime of a catchment in a hypothetical future), for forecasting under the assumption of perfect inputs (where the driving hydrologic forecast is independently obtained via a numerical weather prediction and a hydrologic model) or for regionalization of the WQM (where the model is transferred to a catchment without data). In all of these cases, an ideal model should have unbiased parameter estimates. This is our goal in identifying the optimal input errors, not to use the model for predictions with input data suffering the same errors.

## 2.2 The ~~innovation-introduction~~ of the secant method

Considering the limitations of BATEA and IBUNE framework discussed in the introduction, an improved strategy should be explored to avoid the high dimension challenge and meanwhile promote the error estimation accuracy. ~~The secant method can be applied to address this problem. This is an iterative process to produce better approximations to the roots of a real valued equation (Ralston and Jennrich, 1978). Here, the root is the optimal value of each input error and the equation is the corresponding model residual equal to zero. A traditional approach to updating this is impractical because the estimated input error will fully complement the model error and always lead to a zero residual error regardless of the model parameters. More discussion on this is stated in Sect. 4.1.~~

This study attempts to transform the input error quantification into the rank domain to realize it. Here, the rank is defined as the order of any individual value relative to the other sampled values, and determines the relative magnitude of each error in

all data errors. For example, in the 1<sup>st</sup> iteration in Table A 1, the error at 15<sup>th</sup> time step, -0.29, is the smallest value among all the sampled errors, therefore, its rank is 1. In current methods, an assumption of input error model is necessary to set, which provides an overall distribution for the estimated input errors. If there is knowledge of the error distribution (i.e. cumulative distribution function (CDF) of input errors), the error value only depends on its rank in this distribution. Therefore, under the condition of a certain input error model, the rank estimation will bring similar results as the direct value estimation. Besides, the rank estimation has a few advantages over the direct value estimation. The discussion on this is stated in Sect. 4.1.

In the rank domain, the challenge turns to find a way to effectively adjust the input error rank to minimize the residual error. The secant method can be applied to address this problem. This is an iterative process to produce better approximations to the roots of a real-valued equation (Ralston and Jennrich, 1978). Here, the root is the optimal value of each input error and the equation is the corresponding model residual equal to zero, and this error distribution can then constrain the value range of sampled errors. Therefore, the secant method is very useful in the rank domain, where the root turns to the optimal rank of each input error (rather than its value) and the equation is still the corresponding model residual equal to zero. This new approach, referred to as the Bayesian error analysis with reshuffling (BEAR) method, should be implemented in two steps: sampling the errors from the estimated error distribution and reshuffling these sampled errors corresponding to the inferred error ranks via the secant method.

The secant method (Ralston and Jennrich, 1978) can be repeated as

$$k_{i,q} = k_{i,q-1} - \varepsilon_{i,q-1}^p \frac{k_{i,q-1} - k_{i,q-2}}{\varepsilon_{i,q-1}^p - \varepsilon_{i,q-2}^p} \quad (5)(4)$$

until a sufficiently accurate target value is reached. In this study, the target value is a residual of zero ( $\varepsilon_{i,q}^p = 0$ ) indicating a perfect model fit with input errors estimated exactly. Here,  $k_{i,q}$  represents the estimated rank for  $i$ th input error at the  $q$ th iteration,  $\varepsilon_{i,q-1}^p$  is the residuals corresponding to the input error rank  $k_{i,q-1}$ . The error rank of each data point is updated respectively via Eq. , where  $i = 1, \dots, n$ .  $n$  is the data length and also the number of the estimated errors as these errors are data-based.

After calculating Eq. , it is possible that the rank  $k_{i,q}$  is out of the rank range (for example, less than 1 or more than  $n$ ), or not an integer. Sorting  $k_{i,q}$  in all the ranks  $k_{i,q} (i = 1, \dots, n)$  can address this problem by effectively assigning to each of them a new integer rank based on its position in the sorted list, sealing the calculated ranks  $k_{i,q}$  to an integer from 1 to  $n$ . Thus, in Eq. (5),  $k_{i,q}$  should be changed to  $K_{i,q}$ , representing the pre-rank. After sorting  $K_{i,q}$  for all the errors, the post-rank  $k_{i,q}$  will then belong to reasonable values. The specific calculation of the error rank is demonstrated in the 7th and 8th row in Table A 1.

145 From the above, estimating the rank of input errors via the secant method can be described as the following two equations two steps:

Update the rank of each input error  $K_{i,q}$  ( $i=1,\dots,n$ ) via the secant method respectively for  $i=1,\dots,n$  :

$$K_{i,q} = k_{i,q-1} - \varepsilon_{i,q-1}^p \frac{k_{i,q-1} - k_{i,q-2}}{\varepsilon_{i,q-1}^p - \varepsilon_{i,q-2}^p} \quad (6)(5)$$

Sorting  $K_{i,q}$  ( $i=1,\dots,n$ ) in all the error pre-ranks  $K_q$  to obtain a reasonable rank:

150 
$$k_{i,q} = k(K_{i,q}) \quad (7)(6)$$

where  $k( )$  means calculating its rank.

Thus, the procedure of input error quantification has been developed via the following key steps: 1) Sample the errors from the assumed error distribution to maintain the overall statistical characteristics of the input errors; 2) Update the input error ranks to minimize the model residual via the secant method (Eq. **Error! Reference source not found.** and ); 3) Reorder these sampled errors according to the updated error ranks; 4) Repeat 2) and 3) for a few iterations until a defined target is achieved. This new algorithm is referred to as Bayesian error analysis with reordering (BEAR). An example to illustrate how the BEAR method works is presented in Appendix A.

### **2.3 Approximate Bayesian Computation – Sequential Monte Carlo (ABC – SMC)**

This study chooses Approximate Bayesian Computation via Sequential Monte Carlo (ABC – SMC) as the calibration scheme. ABC – SMC was first proposed by Sisson et al. (2007) and developed in the research of Toni et al. (2008). The ABC method is especially useful for problems in which the likelihood function is analytically intractable or costly to compute in traditional Bayesian approaches. For formal Bayesian approaches,, the likelihood function must be set carefully to meet the assumption about the residual error distribution, and this setting impacts the parameter estimation (Smith et al., 2015, McInerney et al., 2017, Wu et al., 2019). In the ABC method, setting an objective function is more general allowing for potentially complex input error distributions where the likelihood is difficult to write.

In the ABC – SMC approach, the parameter  $\theta^p$  is first sampled from a prior distribution  $P(\theta^p)$  (referred to as population 1).

Then it is propagated through a sequence of intermediate distributions  $P(\theta^p | OF(\mathbf{Y}^o, \mathbf{Y}^p) \leq \tau_s)$ ,  $s=1,\dots, F-1$  (referred to as intermediate population 2, ..., F-1), until it represents a sample from the target distribution  $P(\theta^p | OF(\mathbf{Y}^o, \mathbf{Y}^p) \leq \tau_F)$  (referred

to as the posterior distribution). The tolerance  $\tau_s$  of the objective function is chosen that  $\tau_1 > \dots > \tau_F > 0$ , thus the distributions sequentially evolve towards the target posterior.

## 2.4 Algorithm and an example of the BEAR method

According to the previous derivations, the algorithm quantifying input errors via the BEAR method is demonstrated in Fig. 1 and an illustrative example is presented in Table 1 and Fig. 2. Based on an ABC-SMC calibration scheme, the BEAR method works by replacing the observed input with a modified input that is obtained through the estimated input error rank via the secant method. In Fig. 1,  $s$  refers to the number of the sequential updating populations in the ABC-SMC scheme, which increases until the objective function (measuring the fit between the calibration data and model outputs) of the  $s$ th population is less than the final tolerance  $\bar{\tau}$ . The final tolerance  $\bar{\tau}$  (i.e. the stopping criterion) is difficult to set before calibration due to the unknown range of objective function values, but in practice, it can be estimated after several population calibrations, according to the actual calculation range of the objective function and the target accuracy. In this study, the calibration stops when 1000 proposed parameter sets are rejected in a row. The first tolerance  $\bar{\tau}_1$  should be set sufficiently large to start the update. Any intermediate tolerance  $\bar{\tau}_s$  is set as the 30% quantile of the objective function results of the previous population  $s-1$ , such that it reduces automatically with a new population calculation.

In each calibration population, the input error ranks are updated over  $q$  iterations, where  $q$  increases until the objective function is less than tolerance  $\bar{\tau}_s$ . When  $q=1$  and  $q=2$ , the input errors are randomly sampled from the estimated error distribution because two sets of samples are prerequisites for the updating via the secant method (Table 1). Regarding these, a series of error ranks  $\mathbf{k}_q^p$ , modified inputs  $\mathbf{X}_q^p$ , model outputs  $\mathbf{Y}_q^p$ , and model residuals  $\mathbf{e}_q^p$  are calculated, demonstrated as the 1st and 2nd iteration in Table 1. In later iterations ( $q \geq 3$ ), the error rank  $\mathbf{k}_q^p$  is updated via the secant method (Eq. (6) and (7)), demonstrated in the first two columns in the 3rd and 4th iteration in Table 1. According to the new rank  $\mathbf{k}_{i,q}^p$ , the value with the same rank in the 2nd iteration is the estimated error in the new iteration. For example, the new rank at the 1st time step in the 3rd iteration is 6, and its corresponding value in the 2nd iteration is 0.02, therefore, 0.02 is set as the updated input error at the 1st time step in the 3rd iteration. After the same reshuffling strategy, the re-ranked input errors will then lead to a new series of the modified inputs  $\mathbf{X}_q^p$ , model outputs  $\mathbf{Y}_q^p$  and model residuals  $\mathbf{e}_q^p$ .

Note however if the model parameters are far away from the true values, especially in the initial population, iterative updating of the error ranks will have little effect in reducing the model residual. Therefore, the maximum times number of iterations should be set, referred to as  $Q$ .  $Q$  is set as 20 in this study. If  $q$  exceeds  $Q$ , the algorithm returns to the mutation step resampling the model parameters (seen in in Fig. 1). An example of four iterations is demonstrated in Table 1 and Fig. 2.

In the example given in Table 1, before reshuffling errors (i.e. the 1st iteration and 2nd iteration), the input errors do not approach the true values shown in Fig. 2, having much larger objective function results than the 3rd and 4th iteration. After the error reshuffling, the objective function calculated in the 4th iteration is smaller than the result in the 3rd iteration, illustrating that the estimated errors in the 4th iteration are closer to the true values than the 3rd iteration. This is also supported

by Fig. 2 where the red line (4th iteration) has a stronger correlation with the black line (true input error) than the yellow line (3rd iteration). From the above, the true input errors can be approximated through updating the error ranks to minimize the objective function of the residuals.

### **2.3 Integrating the BEAR method into the Sequential Monte Carlo approach**

205 The core strategy of the BEAR method is to identify the input errors by estimating their ranks, which can be easily integrated into formal Bayesian inference schemes (for example, Markov chain Monte Carlo (MCMC, (Marshall et al., 2004)) and Sequential Monte Carlo (SMC, (Jeremiah et al., 2011, Del Moral et al., 2006))) and other calibration schemes (for example, the generalized likelihood uncertainty estimation (GLUE, (Beven and Binley, 1992))). Based on the traditional calibration approach, the BEAR method works by replacing the observed input with a modified input that is obtained through the estimated  
210 input error rank via the secant method. This study applies the SMC sampler and derives the BEAR method from a Bayesian theoretical foundation in Appendix B. In the SMC approach, the model parameter is first sampled from a prior distribution and then propagated through a sequence of intermediate populations by repeatedly implementing the reweighting, mutation and resampling processes, until the desired posterior distribution is achieved (Del Moral et al., 2006). The details of the SMC algorithm can be found in the study of Jeremiah et al. (2011).

215 **Error! Reference source not found.** demonstrates the integration of the BEAR method into the SMC sampler. In the SMC scheme,  $s$  refers to the number of sequential populations. A population means a group of parameter vectors (particles) that is updated in each iteration. The maximum number of the population  $S$  is set as 200 in this study. In each sequential population,  $N$  particles of model parameters are calibrated.  $N$  is set as 100 in this study. For each particle of the model parameters, the corresponding input error ranks are updated over  $q$  iterations, where  $q$  increases until the acceptance probability is larger than  
220 a number randomly sampled from 0 to 1. It should be noted that if the model parameters are far away from the true values, especially in the initial population, iterative updating of the error ranks will have little effect in reducing the model residual. Therefore, the maximum number of iterations should be set, referred to as  $Q$ .  $Q$  is set as 20 in this study. If  $q$  exceeds  $Q$ , the algorithm returns to the mutation step in Fig. 1.

#### **2.5.2.4 Comparison with other methods**

225 The application of the BATEA framework is limited by high dimension computation (Renard et al., 2009). It probably becomes impractical in quantifying the data-varying errors (rather than the event-varying errors in the study of BATEA (Kavetski et al., 2006)), where the dimension easily exceeds 1000 (Haario et al., 2005). Therefore, the BATEA method is not considered in the comparison. In this study, three methods, including the “Traditional” method, “IBUNE” method and “BEAR” method, are compared to evaluate the ability of the BEAR method in estimating the model parameters and quantifying input errors.  
230 To evaluate the ability of the BEAR method in quantifying input errors, three methods are compared, denoted as method T, D, R. “Traditional” method regards Method “T” is the “traditional” method, regarding the observed input as error-free without



identifying input errors (i.e. Eq. (2)), while the other two methods employ a latent variable to counteract the impacts of input error and build the modified input (i.e. Eq. (4)). (3). In the “IBUNE” method, “D” refers to the probability “Distribution” of input error, which is additional information considered in the calibration. This error distribution can be estimated before calibration according to the studies in the introduction. Especially in the context of proxy errors, the probability distribution can be easily calculated via the residuals between the measurements and the corresponding proxy values. From this error distribution, potential input errors are randomly sampled and filtered by the minimization of the objective function, which is similar to the basic framework of the IBUNE method (Ajami et al., 2007). Although the comprehensive IBUNE framework additionally deals with the model structural uncertainty via the Bayesian Model Averaging (BMA) method, this study only compares the capacity of its input error identification approach. The “BEAR” method adds a reordering process into the “IBUNE” method to improve the accuracy of input error quantification. Method R represents the BEAR method developed in this study. “R” refers to the “Reshuffling” strategy via the secant method, which is an additional process to that used in method D to improve the input error quantification.

### 3 Case studies

#### 3.1 Water quality model: the build-up/wash-off model (BwMod)

This study tests the BEAR algorithm in the context of the build-up/wash-off model (BwMod), which is a group of models to simulate two processes in sediment dynamics, including the build-up of sediments during dry periods and the wash-off process during wet periods. The two formulations were developed in a small-scale experiment (Sartor and Boyd, 1972), while in applications at the catchment scale, the conceptualized parameters largely abandon their physical meanings and the formulations can be considered a “black-box” (Bonhomme and Petrucci, 2017). This study chooses Eq. (4) to describe the build-up process and Eq. (5) to express the wash-off of sediments, representing the non-linear relationship between the wash-off load (output) and the runoff-rate (input). These two equations were applied in the research of Sikorska et al. (2015) and in this study, are written in the MATLAB programming language with the integration of the BEAR method. The time scale is typically set as daily, and the spatial scale is set as the catchment in this study. This version of BwMod has four parameters (Table 2). Sikorska et al. (2015) and in this study are integrated with the BEAR method. This study will test the BEAR algorithm in a case of simulating the daily sediment dynamics of one catchment, thus, the time scale is typically set as daily and the spatial scale is set as the catchment. This version of BwMod has four parameters (Error! Reference source not found.). The model input is streamflow, which typically comes from the observation of a rating curve. As discussed in the introduction, the error distribution can be estimated prior to the model calibration via a rating curve analysis. The output of the BwMod is the concentration of total suspended solids (TSS), whose transport can be efficiently simulated by the conceptualization of the build-up/wash-off process (Bonhomme and Petrucci, 2017, Sikorska et al., 2015). Although BwMod is relatively simple compared with process-based WQMs, its nonlinearity and the use of surrogates for the input data can make it a typical WQM scenario to test the BEAR algorithm.

The overall BwMod equations are:

265

$$\frac{dS_{a,t}}{dt} = \kappa \cdot (S_{max} - S_{a,t}) - s(S_{a,t}) \quad \frac{dS_{a,t}}{dt} = \kappa \cdot (S_{max} - S_{a,t}) - s(S_{a,t}) \quad (4)$$

where the descriptions of  $\kappa$  and  $S_{max}$  are shown in **Error! Reference source not found.**,  $S_{a,t}$  (kg) is the sediment amount available on the catchment surface to be washed-off at time  $t$ ;  $s(S_{a,t})$  (kg/s) is the amount of sediment in the stream at time  $t$ , described by the function

$$s(S_{a,t}) = a \cdot (Q_t)^b \cdot S_{a,t} \quad (5)$$

270

where the descriptions of  $a$  and  $b$  are shown in **Error! Reference source not found.**, and  $Q_t$  is the streamflow at the catchment outlet at time  $t$ .

The output TSS concentration  $C_{TSS,t}$  (kg/m<sup>3</sup>) is derived via:

$$C_{TSS,t} = \frac{s(S_{a,t})}{Q_t} \quad (6)$$

### 3.2 Case study 1: Synthetic data suffering from input errors and parameter errors

275

First, the BEAR method is testedTo test the capability of the secant method in identifying the input error ranks in the process of the model parameter estimation, the BEAR method is first implemented in a controlled situation with synthetic data, where the model is affected only by input errors and parameter errors. The true input  $X^*$  is set as the daily streamflow data of the catchment in the real case (USGS ID: 04087030), covering 1095 days from 2009/10/01 to 2012/09/29. The true output  $Y^*$  is the simulated TSS concentration via BwMod corresponding to the true input  $X^*$  and model parameters set as the reference values in **Error! Reference source not found.** In case study 1, the observed output  $Y^o$  is assumed to be the same as the true simulation  $Y^*$ , i.e. without error. The observed input  $X^o$  is generated based on two types of input error models: an additive formulation and a multiplicative formulation, and the errors are assumed to follow a normal distribution with mean  $\mu$  as 0.2 and standard deviation (SD)  $\sigma$  as 0.5. If the input errors are estimated based on a rating curve, like the procedure in the following real case, the mean of input error should be 0. But in order to test the ability of the BEAR method in wider applications, a systematic bias 0.2 has been considered in the synthetic case even though this is unlikely to manifest in real situations. An additive formulation (denoted as ‘add’ in **Table 3****Error! Reference source not found.**) is suitable to illustrate the error generation in measurements, while the multiplicative formulation (denoted as ‘mul’ in **Table 3****Error! Reference**

285

source not found.) is specifically applied for errors induced from a log-log regression procedure, which is common for water quality proxy processes (Rode and Suhr, 2007). In the additive formulation, the generated input may be negative. If so, the negative input should be truncated to a positive value. In the multiplicative formulation, the generated input will stay positive. Given the description in the introduction, the input error model can be pre-estimated independent of calibration by analysing the input data in some studies. While in other cases, the input error model cannot be estimated or its accuracy is in question. Therefore, two scenarios about the prior information of  $\sigma$  have been considered: one is fixed as the reference values (denoted as 'fixed' in Error! Reference source not found.), the other one is estimated as the hyperparameters with the model parameters (denoted as 'inferred' in Error! Reference source not found.). Therefore, Synthetic case 1 considers four scenarios, including two sets of input data generating from two input error models and two types of prior information about the error parameter  $\sigma$  (the details are shown in Error! Reference source not found.).

Each scenario is calibrated via the traditional method, the IBUNE method and the BEAR method respectively. In the calibration, the objective function is set as the Mean Squared Error (MSE)-2.4. Considering the unknown initial sediment loads in real applications, the calibration sets 90 days as a warm-up period to remove the influence of antecedent conditions. Following the algorithm described in Sect. 2.4, the model parameters and the time-varying input errors are estimated. In each population of the ABC-SMC calibration scheme, 50 sets of model parameters are updated. In the first population, the model parameters are sampled from a uniform distribution with the prior range described in Table 2.

The prior information about error parameters (i.e.  $\sigma$  and  $\mu$ ) contains two conditions: one is fixed as the reference values (denoted as 'fixed' in Table 3), the other one is given the prior range, which needs to infer the error parameters in the calibration (denoted as 'inferred' in Table 3).

To sum up, this study considers four scenarios in the synthetic case, including two sets of synthetic data generating from two input error models and two types of prior information about the error parameter (the details are shown in Table 3). Each scenario is calibrated via method T, method D and method R respectively. Their algorithms are described in Sect. 2.5 and their results are compared in Fig. 3 and Fig. 4. Figure 3 shows the statistical characteristics of the overall estimations. Figure 4. To compare the ability of different methods in estimating the input error and model parameter, this study selects the following statistical characteristics. The SD of the estimated input errors represents the accuracy of the input error distribution (0.5 is the reference value). The correlation between the estimated input error and the true input error evaluates the capability of the method in catching the temporal dynamics of input error. The Nash-Sutcliffe efficiency (NSE) of the modified input vs true input measures the precision of the input data after removing the estimated input errors. In the calibration part, the simulated output corresponds to the modified input and estimated model parameters, and its NSE compared to the true output measures the goodness-of-fit. In the validation part, the simulated output corresponds to the true input and estimated model parameters, and its NSE compared to the true output can assess the accuracy of the model parameter estimation. These statistical characteristics are calculated as the weighted-average values considering the weights of each estimation in the posterior

320 distribution and compared in Fig. 2. Figure C 1 in Appendix C demonstrates the temporal dynamics of input estimations and model simulations-

Evaluating the input error quantification, method R always has much higher correlations with the true error series than method D in all calibration scenarios (shown in Fig. 3(3)). When the error parameters are inferred, the estimations of  $\sigma$  via method D are smaller than the reference value (shown in Fig. 3(1)). This conclusion has also been reported in the study of Renard et al. (2009). The reason for this is that the randomness of the likelihood function leads to an underestimation of the SD of input errors. Compared with method D, the  $\sigma$  estimation via method R is less biased from its true value (shown in Fig. 3(1)), while the estimation of synthetic case 1. In Fig. C1, “Reliability” is the ratio of observations caught by the confidence interval of 2.5%-97.5%  ~~$\mu$~~  is worse via method R (shown in Fig. 3(2)).

Evaluating the model simulation, the BEAR method always produces the best output fit in all scenarios, supported by the highest green bars in Fig. 2(4). Although its correlations with the true error series are much higher than the IBUNE method (red bars) in all scenarios (in Fig. 2(2)), the BEAR method cannot ensure a better input estimation (in Fig. 2(3)) and its ability depends on the prior information of the input error parameter. When the error parameters are fixed at the reference values (in the scenarios *add-fixed* and *mul-fixed*), the BEAR method always outperforms the other two methods in the input modification and model parameter estimation, as its NSE is the highest (green bars in Fig. 2(3) and (5)). Without the reordering strategy, the IBUNE method even gives worse input modification, model simulation and parameter estimation than the traditional method, demonstrated by the lower red bars than blue bars in Fig. 2(3), (4) and (5). Evaluating the model simulation, method R always produces the best output fit in all scenarios, supported by the highest red boxplots in Fig. 3(4). Also in Fig. 4, regardless of the calibration scenarios, the output uncertainty bands of method R (red parts) almost overlaps the true output (green line), being much better than method T (pink parts) and method D (blue parts). However, the input uncertainty bands vary depending on the calibration scenarios. When the error parameters are fixed at the reference values (in the scenarios *add-fixed* and *mul-fixed*), method R always outperforms the other two methods regardless of input error models, as its Nash-Sutcliffe efficiency coefficient (NSE) are the highest (shown in Fig. 3(5)). In Fig. 4(1) and Fig. 4(3), the input uncertainty bands of method R (red parts) generally converge to the true value (green line), being better than method D (blue parts). Without the reshuffling strategy, Method D even gives worse input estimation and model simulations than method T, demonstrated by the lower blue boxplots than pink boxplots in Fig. 3(5)) and Fig. 3(4). This illustrates that the ill-posed error sources in method D exert a negative impact on the model simulations. When the error parameters are inferred (in the scenarios of *add-inferred* and *mul-inferred*), the IBUNE method can improve the input data and the model parameter estimation compared with the traditional method (in Fig. 2(3) and (5)) although the estimations of  $\sigma$  via the IBUNE method are always smaller than the reference value (in Fig. 2(1)). This result has also been reported in the study of Renard et al. (2009), which indicates that the randomness of the likelihood function leads to an underestimation of  $\sigma$  of input errors. Unlike the IBUNE method, the performance of the BEAR method depends on the setting of the input error model. In the *add-inferred* scenario,

the BEAR method is still better than other methods, having a bigger NSE (in Fig. 2(3), (4) and (5)) and the closer  $\sigma$  estimation to reference value (in Fig. 2(1)). While in the *mul-inferred* scenario, the modified inputs and estimated parameters via the BEAR method are worse than the IBUNE method (in Fig. 2(3) and (5)).

355 performance of method R depends on the input error models. For the scenario of *add-inferred*, method R is still better than other methods, having the biggest NSE (shown in Fig. 3(5)) and the closest error parameter estimation to the reference value (shown in Fig. 3(1) and Fig. 3(2)), although the input uncertainty band is more negatively biased from the true value (green line) than method D in Fig. 4(2). For the scenario of *mul-inferred*, the modified inputs via method R are further from the reference value than method D (shown in Fig. 3(5)), which might result from worse  $\mu$  estimations for the input error (shown in Fig. 3(2)).

### 3.3 Case study 2: Synthetic data suffering from input errors, parameter errors and output observation errors

Case study 1 is an ideal situation that is used to test the effectiveness of the BEAR method in isolating the input error and the model parameter error. However, in real-life cases, other sources of errors (i.e. model structural error and output data error) will impact this effectiveness. To explore the ability of the BEAR method with the interference of other sources of errors, the output observational errors with the increasing standard deviations are considered to build the synthetic data based on the scenario 3 and 4 in the case study 1 (the details has been shown in **Error! Reference source not found.**).

**Error! Reference source not found.** demonstrates in the *mul-fixed* scenario where the prior information of standard deviation of input errors is accurate, the BEAR method always brings a better input modification than other methods, although its ability is impaired by the impact of the output observational errors as the NSEs reduces with the increasing SD of the output observational error. The IBUNE method leads to an even worse modified input than the input data without modification in the Traditional method. In the *mul-inferred* scenario where the standard deviation of input errors cannot pre-estimated accurately and given in a wide range, the BEAR method brings worse input data while the IBUNE method can modify the input data.

#### 3.3.4 Case study 23: Real data

To explore the ability of the BEAR method in real-life applications, a real case of one catchment located in southeast Wisconsin, USA is demonstrated. **Error! Reference source not found.** is a description of the test catchment and data (Baldwin et al., 2013). The daily TSS concentration and streamflow data are collected from the USGS database on National Real-Time Water Quality (<https://nrtwq.usgs.gov/>). The daily streamflow data in the USGS database comes from a stage-streamflow rating curve, where the stage and streamflow form a log-log linear relationship and the streamflow proxy errors follow a normal distribution with  $\mu$  as 0 and  $\sigma$  as 0.103. This prior information is used in the real calibration, denoted as *O-fixed* scenario in **Error! Reference source not found.**, where “O” represents the input data that comes from the observations of the rating curve. According to the results of **Error! Reference source not found.** and the assumption of the methodology derivation, the BEAR method is implemented under the assumption that works better when the input

uncertainty is some significant ~~that other sources of uncertainties can be ignored,~~ so another input data source with more significant data uncertainty, a streamflow simulation from a hydrological model, has been considered. This study selects GR4J (Perrin et al., 2003) as the hydrological model and calibrates its parameters with the USGS streamflow data as calibration data. If the USGS streamflow data is regarded as the true input data, the residual error after the model calibration can approximate the data error of GR4J simulation, which follows a normal distribution in log space with  $\mu$  as 0 and  $\sigma$  as 0.764. The BwMod calibration using this input data source and the prior information on data error is denoted as *S-fixed* scenario in **Table 3**. **Error! Reference source not found.** where “*S*” represents the input data that comes from the simulations of GR4J model. To explore the ability of the BEAR method in other situations where the prior information about the input error is not sufficient, two scenarios with a wider range of the error parameters has also been considered, denoted as *O-inferred* and *S-inferred* in **Table 3**. **Error! Reference source not found.** The real case is also calibrated via three methods (i.e. the traditional method ~~T~~, the IBUNE method ~~D~~ and the BEAR method ~~R~~) and adopts the same setting of the calibration algorithm as the synthetic case. Figure 5 uses several statistics to evaluate the calibration scenarios. For all scenarios in Fig. 5(b1), method R always produces a better fit to the output data than method D, consistent with the synthetic case shown in Fig. 3(4). In Fig. 5(b2), “Reliability” here is the ratio of observations caught by the confidence interval of 5%–95%, and the average width of this interval band is referred to as “Sharpness” (Yadav et al., 2007, Smith et al., 2010). In the *S-fixed* and *S-inferred* scenarios with significant input errors, the results of method R show much higher reliability with a larger sharpness. However, in the *O-fixed* scenario with insignificant input errors (i.e.  $\sigma=0.103$ ), the reliability and sharpness of method R are smaller than method D. Fig. 5(a) demonstrates that the  $\sigma$ -estimations vary depending on the calibration methods, but stay almost identical between two data sources. This illustrates that the impacts of other sources of errors significantly impair the error quantification, and their impacts are varied for different methods.

In the real case shown in Fig. 6, method R still produces the best fit to the output and the uncertainty band of the modified input via method D is centered on the observed data. In Fig. 6(c), the uncertainty bands of the modified input are consistent in all scenarios except the *O-fixed* scenario with insignificant input errors (i.e.  $\sigma=0.103$ ). The uncertainty bands are closer to the observed streamflow (green line), even in (c3) and (c4) where the input data comes from the simulated streamflow (black line). According to the results of method T in Fig. 6(a), the simulations corresponding to the observed streamflow (in (a1) and (a2)) catch the dynamics of observed TSS concentration better than the simulations corresponding to the simulated streamflow (in (a3) and (a4)). Here, the observed streamflow from the rating curve should be closer to the true input data, and could be regarded as the reference value. Given that, the modified inputs via method R are more reasonable.

**Error! Reference source not found.** (2) demonstrates the BEAR method always produces a better fit to the output data than the IBUNE method, consistent with the synthetic case shown in Fig. 2(4). In Fig.4(3), except for the *O-fixed* scenario, the results of the BEAR method (in green) show much smaller sharpness than the traditional method (in blue) and the IBUNE method (in red) with almost the same reliability. According to the results of the traditional method in Fig. C2, the simulations from the “*O*” streamflow (in (a1)) catch the dynamics of observed TSS concentration better than the simulations from the “*S*”

streamflow (in (a3)). Thus, compared with the simulated streamflow via GR4J (“S” streamflow), the observed streamflow from the rating curve (“O” streamflow) should be closer to the true input data. In Fig. C2, the modified inputs via the BEAR method are closer to the “O” streamflow (blue dots) than the “S” streamflow (pink dots), even in (c3) and (c4) where the original input data comes from the “S” streamflow. However, the modified input via the IBUNE method is always centred on the original input data it uses. Given being always closer to the “O” streamflow, the modified inputs via the BEAR method are more reasonable than the IBUNE method.

## 4 Discussion

### 4.1 The effectiveness of rank estimation

The novelty of the BEAR method lies in transforming a direct error value estimation to an error rank estimation. In a continuous sequence of data, the potential error values have an infinite number of combinations, while the error rank has limited combinations, dependent on the data length. ~~It is far more efficient to estimate the error rank than estimate the error value. Compared with the IBUNE framework (Ajami et al., 2007), the BEAR method additionally infers the error ranks to adjust the order of the sampled errors and reduce their randomness, which significantly improves the accuracy of the error estimation (as demonstrated by much higher NSEs than method D in Fig. 3). The application of the secant method plays an essential role in this by inferring each error rank according to the residual error. For example, in Table A1, the estimated error at the 1<sup>st</sup> time step could be any value. Even under a constrain of the range from the minimized to the maximized sampled errors (i.e. [-0.29,0.16] in the 1<sup>st</sup> iteration), its value estimation still has infinite possibilities due to the continuous nature of the error. In contrast, the rank is discrete, having only 20 possibilities (i.e. the integrity in [1,20]). From this point of view, it is more efficient to estimate the error rank than estimate the error value.~~

~~One thing to note in the rank estimation~~ However, the rank estimation will suffer from the sampling bias problem. The sampling bias problem is that even corresponding to the same rank, the error sampled at different times could be largely different, especially for a small sample size (depending on the data length) or a large  $\sigma$  of the assumed error distribution. This problem can be addressed by selecting the optimal solution from multiple sampling according to the maximum of likelihood function. In three cases of this study, the sample size is larger than 1000, where the sampling bias problem can be neglected and one error sampling is enough. But in some cases where the sample size is small (i.e. around 10), multiple sampling should be taken. ~~The secant method is a successive approximation algorithm and one single iteration cannot guarantee the optimal results. Considering these two points, the BEAR method set  $q$  iterations in the algorithm (Fig. 1).  $q$  increasing until the objective function becomes smaller than the tolerance.~~

~~Note that~~ Besides, to avoid the high-dimension calculation, modifying each input error according to its corresponding residual error only works in the rank domain. In the value domain, if there is no constraint on the estimated input errors, they will fully compensate for the residual error to maximize the likelihood function and subsequently be overfitted. There are two ways to

impose restrictions. One is to regard errors and model parameters as a whole in calibration, like the BATEA framework (Kavetski et al., 2006), resulting in a high dimensional computation. The other is to sample error randomly from the assumed error model, like the IBUNE framework (Ajami et al., 2007), whose precision cannot be guaranteed due to the error randomness. However, in the BEAR method, the inference focuses on the error rank where the value range of the sampled errors can be effectively limited by the assumed error model. Additionally, adjusting the order of the sampled errors according to the inferred error rank can reduce the randomness in the IBUNE framework (Ajami et al., 2007), which significantly improves the accuracy of the error estimation (as demonstrated by much higher NSEs than the IBUNE method in Fig. 2). The reordering step is implemented when the model parameter has been updated and aims to find the optimal input error series corresponding to the minimized residual error. After the reordering step, the optimal input error is a deterministic function of the model parameter. Thus, unlike formal Bayesian inference, the BEAR method does not update the posterior distribution of the input errors, but identifies the input error through the deterministic relationship between the input error and model parameter.

#### 4.2 The impacts of prior information of input error model

Method D employs the same framework as IBUNE (Ajami et al., 2007), taking advantage of stochastic error samples to modify the input observations. In Fig. 4 and Fig. 6, the uncertainty bands of modified inputs (blue parts) encompass the input observations (black line), illustrating that the intrinsic quality of the input observation determines the algorithm performance. Figure 6 demonstrates that if the input error is insignificant in the residual, like in the *O-fixed* and *O-inferred* scenarios of the real case, the resultant simulations will fit the observed output (green line) well. Otherwise, the simulations are far away from the observed outputs (black line) due to inaccurate input observations (in the *S-fixed* and *S-inferred* scenarios in the real case). As per the finding in the previous study of Renard et al. (2010), if the SD of input errors is inferred with the model parameters, method D will underestimate the SD (Fig. 3(1) and Fig. 5(a2)). If the intrinsic SD of input errors is large, a fixed SD cannot improve the input modification and model simulation, demonstrated by a wider band in Fig. 6(b3) than in Fig. 6(b4). If the SD of input errors is small, the prior information will constrain the impacts of other sources of errors. From the above, the data quality is more important than the availability of prior information for method D, especially when the intrinsic SD of the input error is large.

However, the findings in method R are quite different. Although method R infers the input error by minimizing the model residual error, it is much more effective than method D to minimise the residual errors. For the synthetic case (Fig. 4(c)) and real case (Fig. 6(c)), the model simulations via method R (red parts) are very close to the output observations (green line). In other words, the estimated input error mainly depends on the output observations. Therefore, in the real case with the same output observation (Fig. 6(c)), the modified inputs are consistent among the scenarios. Given this, the model structure error plays an important role in estimating the input error.

To constrain the impacts of the other sources of error, accurate prior information about the input error model is important in method R. In the synthetic case, fixed scenarios always produce a higher NSE of the modified input (Fig. 3(5)) and a larger



480 correlation in the estimation error (Fig. 3(5)) than inferred scenarios. This illustrates that prior information can limit the impacts  
of model parameter error. In Fig. 6(a1), the modified inputs in the real case are around the reference value (green line), while  
in Fig. 6(a2), the modified inputs are biased from the reference value (green line). It should be noted that this difference is  
obvious in the scenarios with insignificant input error (where the model structural error is relatively large). When the input  
error is dominant, like the *S-inferred* scenario, method R becomes more effective to estimate the input error, bringing a more  
485 precise estimation of the error SD than the *O-inferred* scenario and similar results to the *S-fixed* scenario.

To sum up, for method R, an accurate input error model can constrain the adverse impacts of the other sources of errors,  
especially when the other sources of error are dominant. But for method D, the input data quality is more important than this  
prior information.

The IBUNE method takes advantage of stochastic error samples to modify the input observations (Ajami et al., 2007). Figure  
490 C demonstrates compared with *O-fixed* and *O-inferred* scenarios, *S-fixed* and *S-inferred* scenarios uses simulated streamflow  
whose input error is more significant, and the resultant simulations (black line) via the IBUNE method are further away from  
the observed outputs (red dots). As per the findings in the previous study of Renard et al. (2010), if the  $\sigma$  of input errors is  
inferred with the model parameters, the IBUNE method will underestimate  $\sigma$  (in Fig. 2(1) and Fig. 4(1)). If  $\sigma$  is fixed via  
prior information, the input modification and model simulation cannot be improved, especially in the scenarios with large  
495 intrinsic  $\sigma$  of input errors (in Fig. 2 and Fig. 3). From the above, the ability of the IBUNE method depends on the input data  
quality and the improvement of the input data and model simulation only happens when the  $\sigma$  of the estimated input error is  
small. The availability of prior information is insignificant for the IBUNE method, especially when the intrinsic  $\sigma$  of the  
input error is large.

However, the findings in the BEAR method are quite different. Accurate prior information about the input error model is  
500 important in the BEAR method. Figure 3 demonstrates *fixed* scenarios calibrated via the BEAR method always produce a  
higher NSE of the modified input than *inferred* scenarios. This is likely because the prior information can constrain the input  
error distribution and reduce the impacts of other sources of errors. The availability of prior information of the input error  
relies on studies about benchmarking observational errors of water quality and hydrologic data, and the selection of a proper  
input error model is important. Comparing the results in Figure 2, when the input error model is an additive formulation, the  
505 BEAR method consistently brings the best performance regardless of the prior information of the error  $\sigma$ . When the input  
error model is a multiplicative formulation, the BEAR method cannot improve the input data if the prior information of the  
error  $\sigma$  is not accurate. This illustrates that the compensating effect between the input error and parameter error is weaker in  
the additive form of the input error. This is probably related to the specific model structure, as the exponent parameter  $b$  in  
*BwMod* has a stronger interaction with the multiplicative errors than the additive errors. Thus, more comprehensive  
510 comparisons should be undertaken to explore the capacity of different input error models in different model applications.

To sum up, the ability of the BEAR method depends on the accuracy of prior information of the input error parameter and the selection of the input error model. The IBUNE method can modify the input data when the standard deviation of the estimated input error is much smaller than the true value. It is most likely to make use of the stochastic errors to improve the original input data, but not effectively identify the input error.

### 515 4.3 The extension to other modeling scenarios

In this study, the BEAR method was developed in the calibration of BwMod at the daily time scale, whose input and output ~~correspond~~can be regarded as the correspondence at each time step. Therefore, in Eq. ~~(6)~~**Error! Reference source not found.**

, the model residual  $\varepsilon_{i,q-1}^p$  and input error rank  $k_{i,q-1}$  are at the same time step  $i$ . If the water quality system exhibits delayed response, the time lag between the forcing data and the response (described as *lag*) should be considered in the algorithm and

520 Eq. ~~(6)~~**Error! Reference source not found.** needs to be modified as Eq. ~~(11)~~**Error! Reference source not found.**

$$K_{i,q} = k_{i,q-1} - \varepsilon_{i+lag,q-1}^p \frac{k_{i,q-1} - k_{i,q-2}}{\varepsilon_{i+lag,q-1}^p - \varepsilon_{i+lag,q-2}^p} \quad (10)(11)$$

If the response caused by an input is not instantaneous but exhibits persistence (i.e. occurs over several time steps), the autocorrelation in the output should be addressed to ensure the independence assumption of the rank updating is satisfied.

Current ways to deal with this problem in hydrologic modelling can provide a reference to the potential modification of the BEAR method. Autocorrelation in the residual errors can be represented by an autoregressive moving average (ARMA) model (Kuczera, 1983) or autoregressive (AR) (Schaeffli et al., 2007, Bates and Campbell, 2001). ~~However, the ability of these~~

525 BEAR method. Autocorrelation in the residual errors can be represented by an autoregressive moving average (ARMA) model (Kuczera, 1983) or autoregressive (AR) (Schaeffli et al., 2007, Bates and Campbell, 2001). ~~However, the ability of these~~

~~approaches needs further discussion in systems with correlated responses. The correlated part of the error is removed from the residual error and the remaining part will be only impacted by the input error. Thus, the correspondence between the input error and the residual error part is ensured and the latter process will be the same as the application of the BEAR method in~~  
 530 ~~this study. Following this idea, the autoregressive (AR) model has been integrated with the BEAR method in the study of Wu et al. (2021) to deal with the autocorrelation of residual errors in a hydrologic model. The results prove this integration is effective to improve the input error estimation.~~

~~However, this treatment may not guarantee the improvement of the input error estimation in this study where the sediment concentrate is simulated at the daily time scale (Figure D 1). At this time scale, one input (streamflow) may not impact the response (sediment concentration) for multiple time steps and autocorrelation may not be well represented via a simple autocorrelation function. When the temporal resolution of the data is high (i.e. minute) and one model output is affected by many inputs, the memory effect may be addressed effectively via the AR model. Therefore, the specific representation of the autocorrelation in the residual error needs further discussion through comparisons in different time scales or with different characteristics in the memory effect.~~

## 540 5 Conclusion

Taking advantage of the prior information of an input error model, a new method, Bayesian error analysis with ~~reshufflingreordering~~ (BEAR), is ~~developedproposed~~ to approach the time-varying input errors in WQM inference. It contains two main processes: sampling the errors from an assumed error distribution and reordering them with the inferred ranks via the secant method. Through the investigation of synthetic data and real data, this method is shown to be ~~robust and effective~~.~~effective but its ability is limited by the accuracy and selection of the input error model.~~ The novelties of this algorithm are: (1) ~~Estimating~~The estimation focuses on the error rank rather than ~~directly estimating~~ the error value, which ~~significantly improves the effectiveness of input error quantification by reducing the potential search space for input errors~~.enhances the constraints of the input error model on the estimated errors and avoids the high dimensionality problem resulting from calibrating all the errors along with the model parameter as a whole. (2) The ~~modification~~introduction of the secant method ~~links~~realizes updating the error rank of each input data according to its corresponding residual,~~which addresses the high dimensionality problem in current calibration methods and tackles the nonlinearity challenge in the WQM transformation.~~

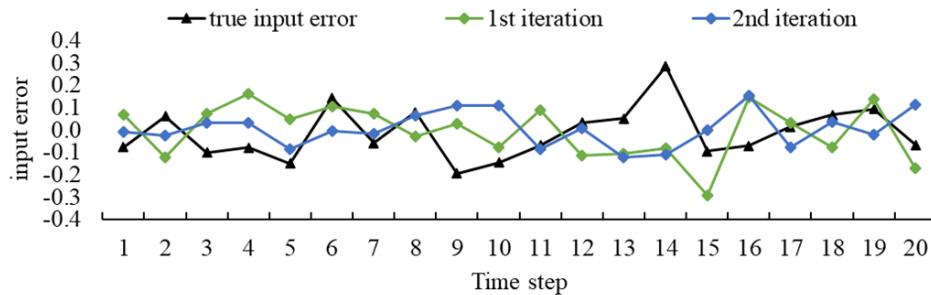
However, the work in this study still identifies a few areas needing to be explored. Firstly, the availability of prior knowledge of the input error model is important. When this information is not reliable or even cannot be estimated, a significant issue is the selection of a suitable error assumption. Thus, a general measure should be found to judge whether an error model is appropriate, especially in real cases where the “true” information is limited. Secondly, extensions of the BEAR method to other water quality modeling scenarios are subject to problems such as delayed and autocorrelated responses. Related studies in hydrologic modeling to deal with the delay and persistency of responses could be references in the modification of the BEAR method. Thirdly, if the sampling and ~~reshufflingreordering~~ strategy is developed within a more comprehensive framework to quantify multiple sources of error, the interactions amongst these error sources might be well-identified and the quantification of individual errors might be improved. This study provides a starting point for developing the rank estimation via the secant method to identify input error. Further study is necessary to modify the algorithm and improve confidence in extended case studies or model scenarios.

## Appendix A: The illustration of the BEAR method

### Table A 1 An example illustrating the BEAR method

		1st iteration (the input errors are randomly sampled)																			
row	time step	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	sampled input error	0.07	-0.12	0.07	0.16	0.05	0.07	0.07	-0.03	0.03	-0.08	0.09	-0.11	-0.11	-0.08	-0.29	0.14	0.03	-0.08	0.14	-0.17
2	input error rank	13	3	14	20	12	17	15	9	10	7	16	4	5	6	1	19	11	8	18	2
3	model residual error	-0.29	0.49	-0.58	-0.98	-0.78	0.29	-0.66	0.59	-1.31	-0.31	-0.87	0.76	0.46	0.54	0.25	-0.80	-0.07	0.56	-0.23	0.40
MSE		0.40																			
		2nd iteration (the input errors are randomly sampled)																			
4	sampled input error	-0.01	-0.02	0.03	0.03	-0.09	0.00	-0.02	0.06	0.11	0.11	-0.09	0.01	-0.12	-0.11	0.00	0.15	-0.08	0.04	-0.02	0.11
5	input error rank	9	14	13	3	10	8	16	17	18	4	12	1	2	11	20	5	15	7	19	
6	model residual error	-0.13	0.23	-0.43	-0.41	-0.21	0.70	-0.23	0.09	-1.88	-1.52	0.20	0.17	0.53	0.60	-0.43	-0.72	0.36	0.12	0.47	-0.82
MSE		0.47																			
		3rd iteration (the error ranks are updated via the secant method)																			
7	calculated pre-rank	5.8	8.7	14.0	8.0	-0.3	22.0	4.3	17.3	-6.1	4.2	6.2	14.3	31.3	42.0	4.7	29.0	10.0	16.9	14.4	7.6
8	ranked rank (post rank)	6	10	12	9	2	17	4	16	1	3	7	14	19	20	5	18	11	15	13	8
		3rd iteration (the input errors are reordered with the updated error ranks)																			
9	reordered input error	-0.02	0.00	0.01	-0.01	-0.11	0.11	-0.09	0.06	-0.12	-0.09	-0.02	0.03	0.11	0.15	-0.08	0.11	0.00	0.04	0.03	-0.02
10	model residual error	-0.23	0.20	-0.34	-0.24	-0.12	0.19	0.14	0.08	-0.40	-0.31	-0.22	0.03	-0.17	0.26	-0.09	-0.55	0.11	0.14	0.27	-0.23
11	MSE	0.06																			

(a) at the 1st and 2nd iteration where the input errors are randomly sampled



(b) at the 3rd iteration where the input errors are reordered according to the updated error ranks

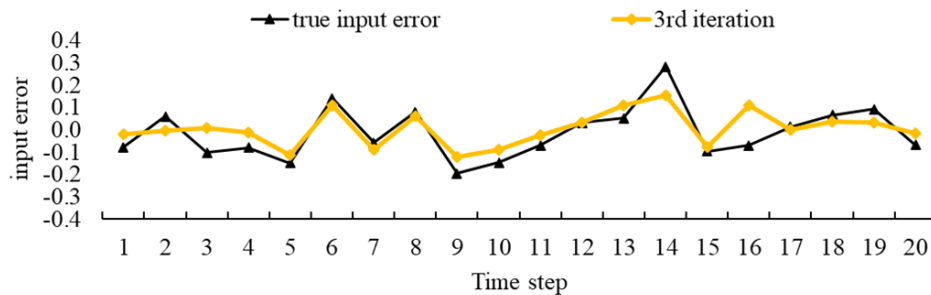


Figure A 1 Demonstration of the input error estimated in Table A.1

The BEAR method for identifying the input errors is implemented after generating the model parameters and contains two main parts: sampling the errors from an assumed error distribution and reordering them with the inferred ranks via the secant method. An example is illustrated in Table A 1 and the explanation about the specific steps is presented in the following contents.

(1) In the 1st iteration ( $q=1$ ), the errors are randomly sampled from the assumed error distribution (row 1), and then are sorted to get their ranks (row 2). This error series is employed to modify the input data, which leads to a new model simulation and model residual (row 3).

(2) Repeat step (1) in the 2nd iteration ( $q=2$ ) as two sets of samples are prerequisites for the updating via the secant method. The results are shown in row 4, 5 and 6. Figure A 1(a) demonstrates that the ranges of the error distribution are the same between the true input errors (black line) and the sampled errors (blue and green lines) as they come from the same error distribution under the condition that prior knowledge of the input error distribution is correct. However, the values at each time step cannot match due to the randomness of the sampling.

(3) At the 1st time step in the 3rd iteration ( $i=1, q=3$  in Eq. **Error! Reference source not found.**), the pre-rank  $K_{1,3}$  is calculated via the secant method (illustrated as the following Eq. **Error! Reference source not found.**). The details are demonstrated in solid boxes in Table A.1.

$$K_{1,3} = k_{1,2} - \varepsilon_{1,2}^p \frac{k_{1,2} - k_{1,1}}{\varepsilon_{1,2}^p - \varepsilon_{1,1}^p} = 9 - (-0.13) \frac{9 - 13}{-0.13 - (-0.29)} = 5.8$$

(4) Repeat step (3) for all the time steps. The calculated pre-ranks are shown in row 7.

(5) Sort all the pre-ranks to get the integral error rank (row 8).

(6) According to the updated error ranks (row 8), the sampled errors in the 2nd iteration (row 4) are reordered. The example for the 1st time step is demonstrated in dotted boxes in Table A.1. The error rank at the 1st time step is updated as 6, and the rank 6 corresponds to the error value -0.02 in the 2nd iteration. Therefore, -0.02 is the input error at the 1st time step in the 3rd iteration. Following this example, the sampled errors at all the time steps are reordered. The results are shown in row 9. Figure A 1(b) demonstrates that after reordering the errors with the inferred ranks, the estimated errors are much closer to the true input error, and the mean square error (MSE) of the model residual reduces in Table A 1.

(7) The reordered input error will lead to a new input data, a new model simulation and a new model residual. The residual result and its MSE statistic are shown in row 10 and 11 respectively.

(8) Check the convergence: If the objective function or likelihood function meets the convergence criterion, stop and the input error estimation is accepted. Otherwise,  $q=q+1$ , repeat step (3)~(8) until  $q$  is larger than the maximum numbers of iteration  $Q$ .

## Appendix B: Theoretical foundation of the BEAR method

### (1) Basic notation

605 In general, a model  $M()$  simulates the output  $Y^s$  given the observed input  $X^o$  and model parameters  $\theta$ , as follows:

$$Y^s = M(X^o, \theta) \quad (1)$$

Here and in the following,  $^s$  represents the simulated value,  $^o$  represents the observed value, and  $^*$  represents the true value.

### (2) Input errors

610 The input errors  $\epsilon_X$  are assumed to be represented by input multipliers, which are sampled from an uncorrelated lognormal distribution, and the observed input  $X^o$  can then be related to the true input  $X^*$  by the following equation:

$$X^o = X^* \exp(\epsilon_X), \epsilon_X \sim N(\mu_X, \sigma_X^2) \quad (2)$$

where  $\epsilon_X$  are assumed to follow a Gaussian distribution with mean  $\mu_X$  and variance  $\sigma_X^2$ .

### (3) Output observational errors and model structural errors

In the derivation, these two parts are assumed to be error-free, therefore,

615 
$$Y^o = Y^* \quad (3)$$

$$M() = M^*() \quad (4)$$

### (4) Remnant errors

Based on the previous assumptions, the observed output equals the true output, and the difference between the simulated output and the observed output,  $\epsilon$ , will be equal to the difference between the simulated output and the true output, as follows:

620 
$$Y^s = Y^o + \epsilon = Y^* + \epsilon, \epsilon \sim (0, \sigma^2) \quad (5)$$

where the remnant errors  $\epsilon$  are assumed to follow a Gaussian distribution with mean 0 and variance  $\sigma^2$ .

### (5) Bayesian inference

According to the study of Renard et al. (2010), the posterior distribution of all inferred quantities is given by Bayes' theorem, as follows:

625

$$\frac{p(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_X, \mu_X, \sigma_X, \sigma | Y^o, X^o) \propto p(Y^o | \boldsymbol{\theta}, \boldsymbol{\varepsilon}_X, X^o) p(\boldsymbol{\varepsilon}_X | \mu_X, \sigma_X) p(\boldsymbol{\theta}, \mu_X, \sigma_X, \sigma)}{p(Y^o | \boldsymbol{\theta}, \boldsymbol{\varepsilon}_X, X^o) p(\boldsymbol{\varepsilon}_X | \mu_X, \sigma_X) p(\boldsymbol{\theta}, \mu_X, \sigma_X, \sigma)} \quad (6)$$

The full posterior distribution comprises the following three parts: the likelihood of the observed output  $p(Y^o | \boldsymbol{\theta}, \boldsymbol{\varepsilon}_X, X^o)$ , the hierarchical parts of the input multiplier  $p(\boldsymbol{\varepsilon}_X | \mu_X, \sigma_X)$  and the prior distribution of deterministic parameters and hyperparameters  $p(\boldsymbol{\theta}, \mu_X, \sigma_X, \sigma)$ .

630

Renard et al. (2009) argue that in the IBUNE method,  $\boldsymbol{\varepsilon}_X$  are randomly sampled in each evaluation of the likelihood function and their different values at different evaluations will lead to the nondeterministic nature of the likelihood function (Equation (6)). In Bayesian inference, the likelihood function should return a fixed value for a given set of arguments. However, the randomness of the likelihood function in the IBUNE method breaks this theoretical foundation. Conversely, in the BEAR method, the secant method is applied to find a deterministic relationship between the rank of each input error and its corresponding model residual error. The residual errors depend on the model parameters  $\boldsymbol{\theta}$ . The magnitude of the whole input errors (i.e. their cumulative distribution function (CDF)) is related to the hyperparameters of the multipliers  $\mu_X, \sigma_X$ . Given the value of each input error is determined by the CDF of the whole input errors and its relative rank among them,  $\boldsymbol{\varepsilon}_X$  depends on  $\mu_X, \sigma_X$  and  $\boldsymbol{\theta}$ , as follows:

635

$$\boldsymbol{\varepsilon}_X = f(\boldsymbol{\theta}, \mu_X, \sigma_X) \quad (7)$$

640

Considering  $\boldsymbol{\varepsilon}_X$  are sampled from  $N(\mu_X, \sigma_X^2)$ ,  $p(\boldsymbol{\varepsilon}_X | \mu_X, \sigma_X)$  is fixed when  $\mu_X, \sigma_X$  are determined and do not need to be considered in Equation (6). Therefore, the posterior distribution of all inferred parameters (Equation (6)) in the BEAR method will turn into:

$$\frac{p(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_X, \mu_X, \sigma_X, \sigma | Y^o, X^o) \propto p(Y^o | \boldsymbol{\theta}, \mu_X, \sigma_X, X^o) p(\boldsymbol{\theta}, \mu_X, \sigma_X, \sigma)}{p(Y^o | \boldsymbol{\theta}, \mu_X, \sigma_X, X^o) p(\boldsymbol{\theta}, \mu_X, \sigma_X, \sigma)} \quad (8)$$

645

The above derivation states if the relationship between the input errors and model parameters (Equation (7)) can be determined, the problem of parameter estimation and input error identification (Equation (6)) can then be interpreted as the updating  $\boldsymbol{\theta}, \mu_X, \sigma_X$  (Equation (8)) in the Bayesian inference. There are two ways to realize this determined relationship: one is to estimate the parameters and input errors together, as the BATEA approach, which will suffer from the high-dimensionality problem (Renard et al., 2010); the other one is to explore the relationship between each input error rank and model parameters

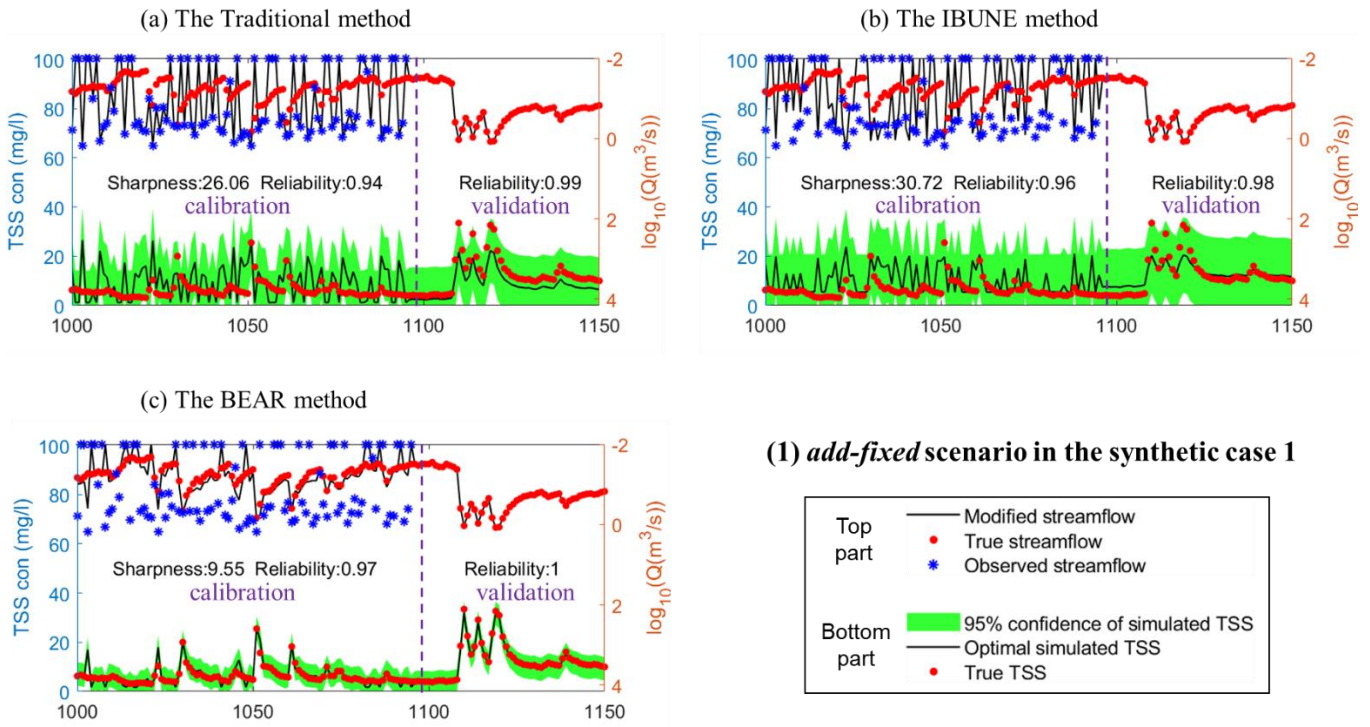


via the secant method first, and then transform the error rank into the error value according to the estimated error parameters

$\mu_x, \sigma_x$ , as the BEAR approach in this study.

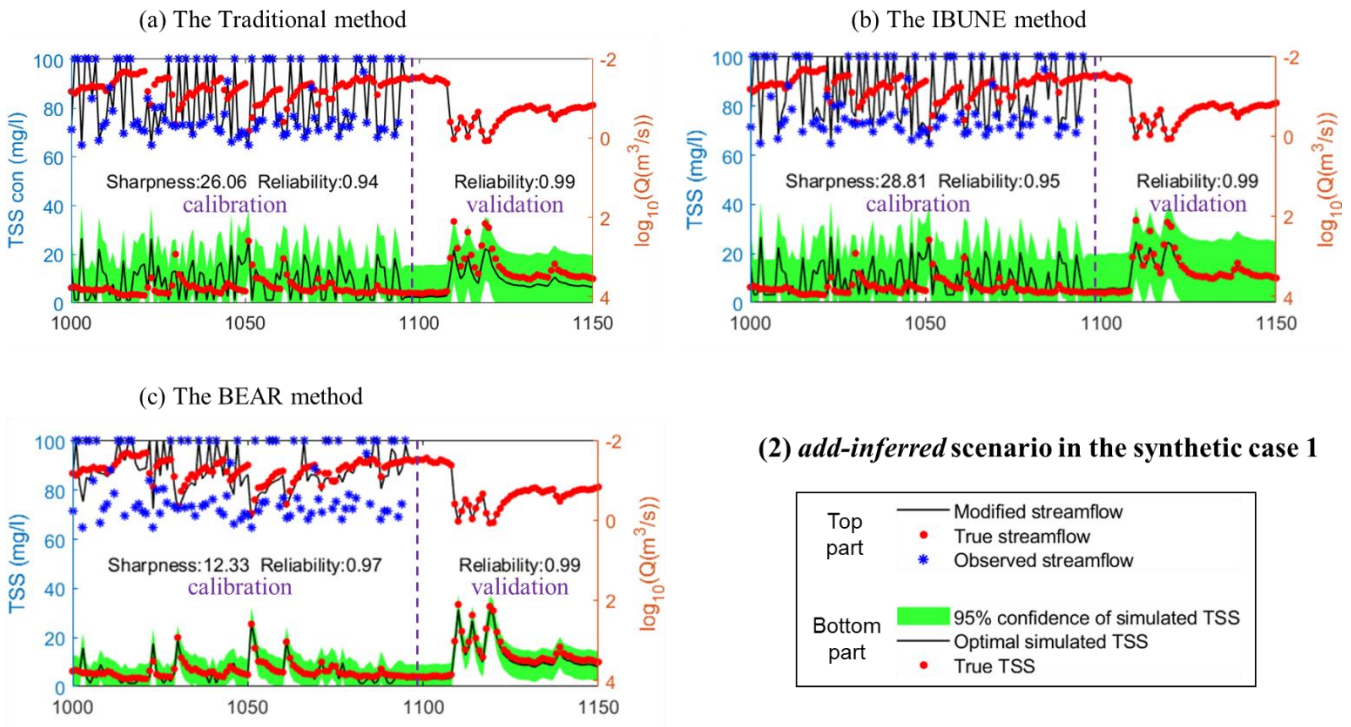
650

**Appendix C: The time series of results in the case study**



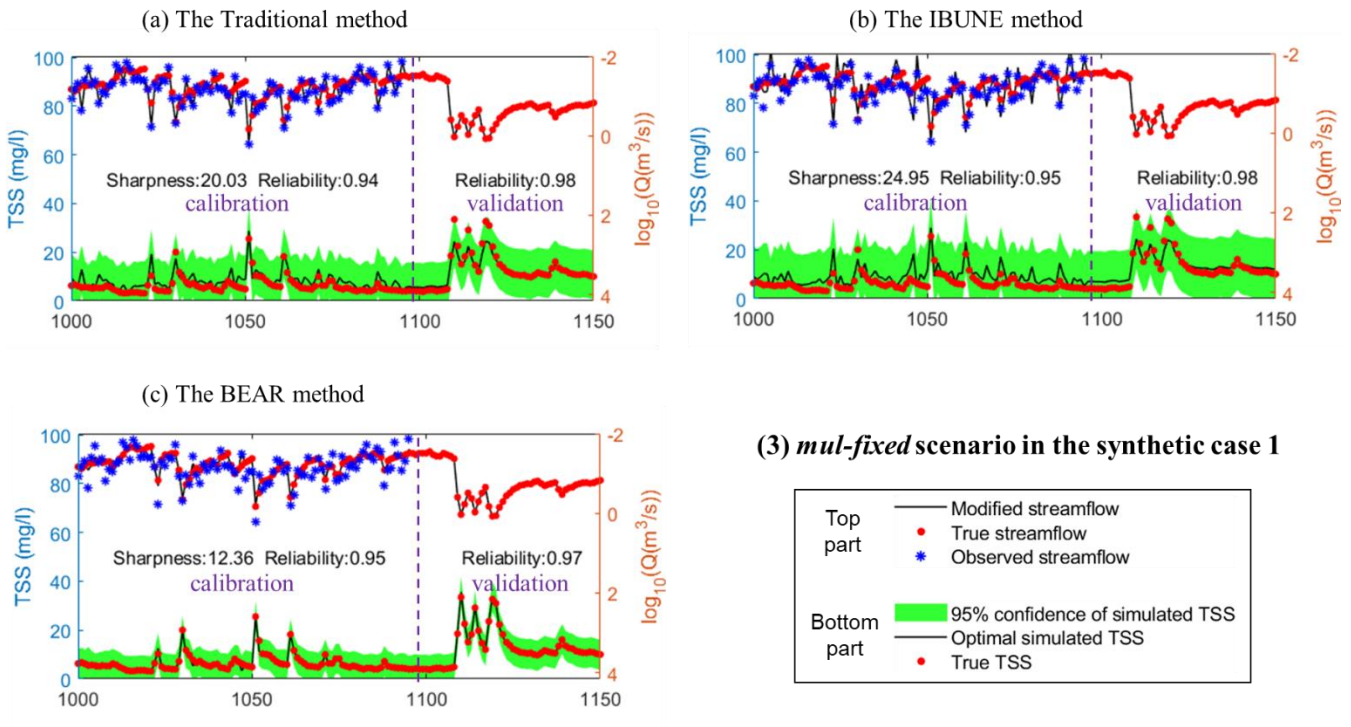
655

**Figure C 1(1) Comparison of time series of synthetic data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method; algorithms are explained in Sect. 2.4) for a select period of *add-fixed* scenarios in the synthetic case 1(notations are given in Error! Reference source not found.)**



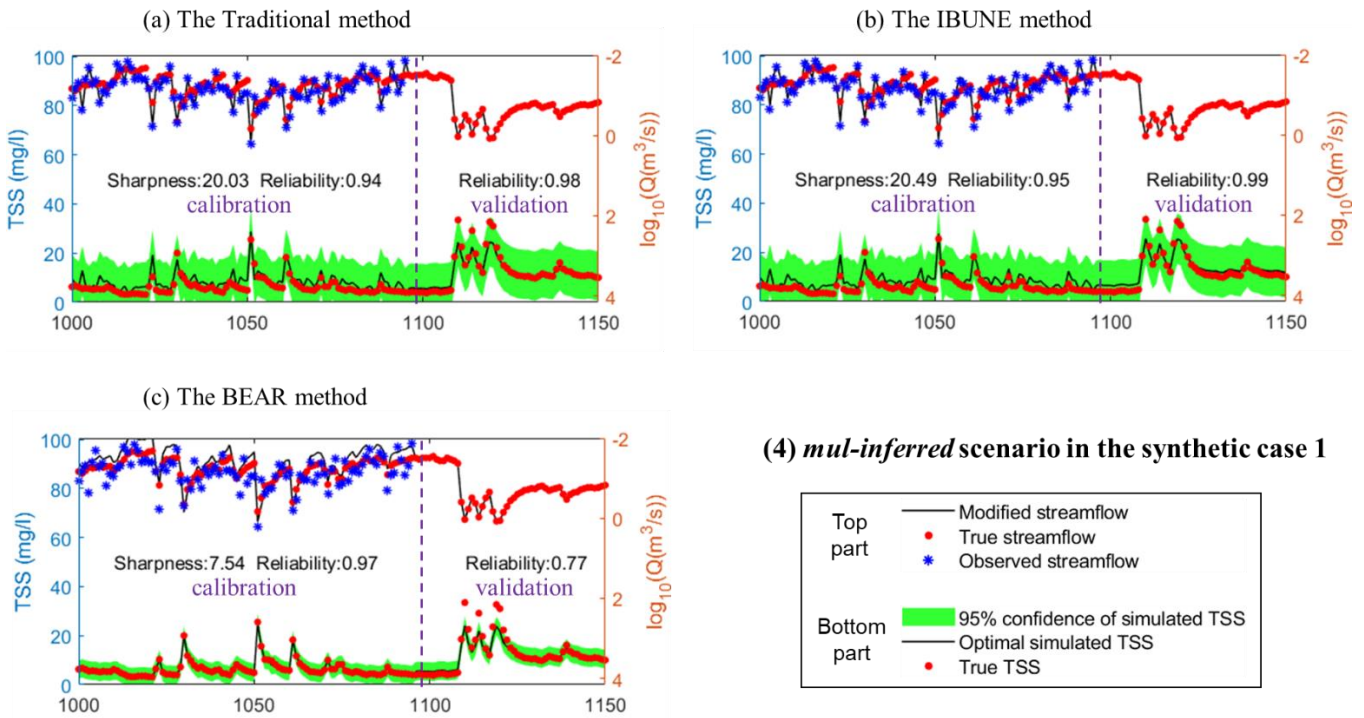
**Figure C 1(2) Comparison of time series of synthetic data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method; algorithms are explained in Sect. 2.4) for a select period of *add-inferred* scenarios in the synthetic case 1(notations are given in Error! Reference source not found.)**

660



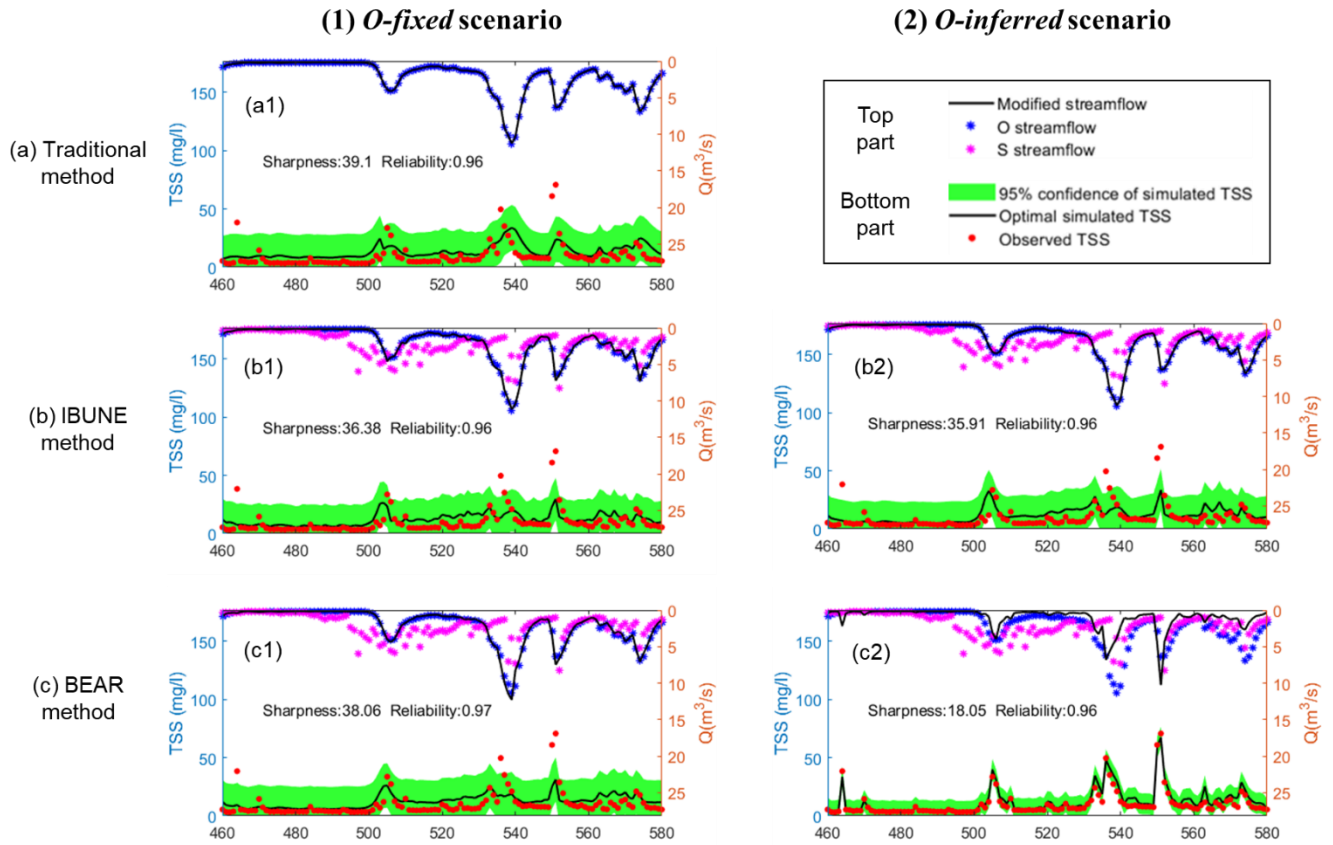
**Figure C 1(3) Comparison of time series of synthetic data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method; algorithms are explained in Sect. 2.4) for a select period of *mul-fixed* scenarios in the synthetic case 1 (notations are given in Error! Reference source not found.)**

665



670

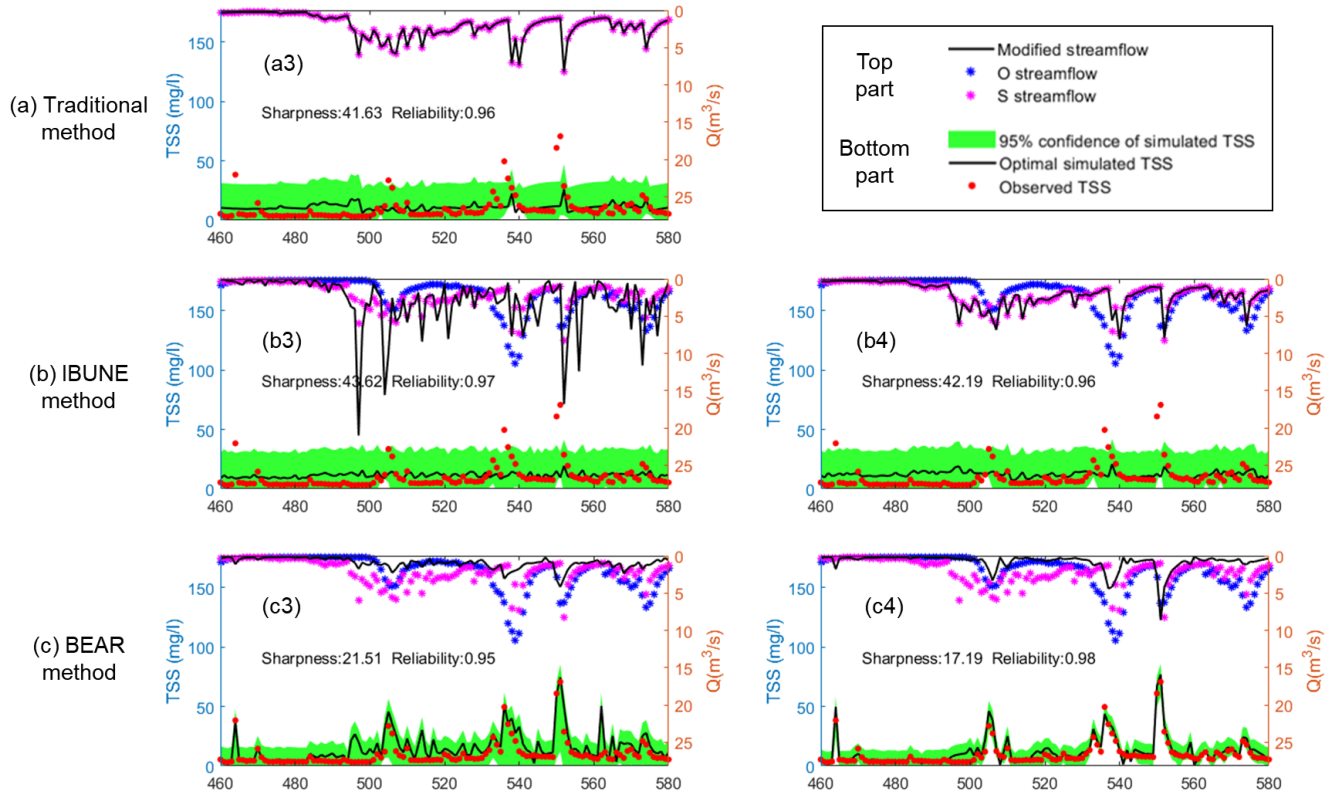
**Figure C 1(4) Comparison of time series of synthetic data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method; algorithms are explained in Sect. 2.4) for a select period of *mul-inferred* scenarios in the synthetic case 1(notations are given in Error! Reference source not found.)**



**Figure C 2(1) Comparison of time series of real data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method, algorithms are explained in Sect. 2.4) for a select period of *O*-fixed, *O*-inferred scenarios in the real case (notations are given in Table 2)**

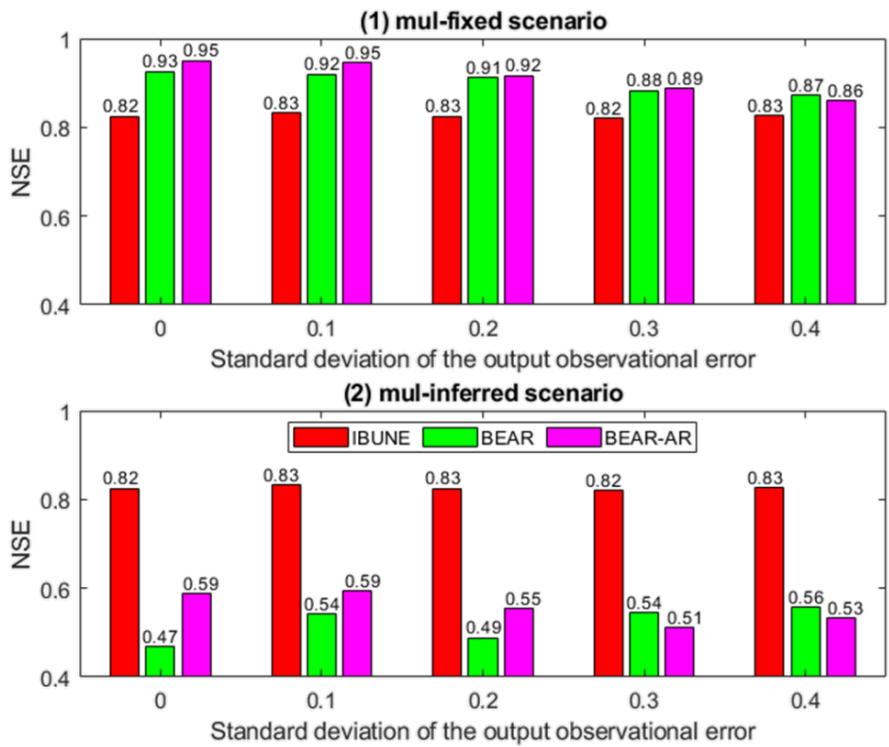
(3) *S-fixed* scenario

(4) *S-inferred* scenario



680

**Figure C 2(2)** Comparison of time series of real data and uncertainty bands estimated via three calibration methods (including the traditional method, the IBUNE method and the BEAR method, algorithms are explained in Sect. 2.4) for a select period of *S-fixed* and *S-inferred* scenarios in the real case (notations are given in Error! Reference source not found.)



690 **Figure D 1 Comparison of Nash-Sutcliffe efficiency (NSE) of the modified input v.s true input under the interference of the output observational errors with the increasing standard deviations in two calibration scenarios in synthetic case 2 (including *mul-fixed* and *mul-inferred*; notations are given in *Error! Reference source not found.*) via three calibration methods (including the IBUNE method and the BEAR method and the BEAR-AR method, the BEAR-AR method is the BEAR method after applying the autoregressive (AR) model to deal with the residual error)**

Table



## Code/Data availability

The daily streamflow and TSS concentration data for real case catchment (ID: USGS 04087030) can be accessed by the  
695 National Real-Time Water Quality website of USGS, the link is <https://nrtwq.usgs.gov/>.

## Author contribution

Lucy Marshall and Ashish Sharma designed the research. Xia Wu developed the research code, analyzed the results, and prepared the manuscript with contributions from all co-authors.

## Competing interests

700 The authors declare that they have no conflict of interest.

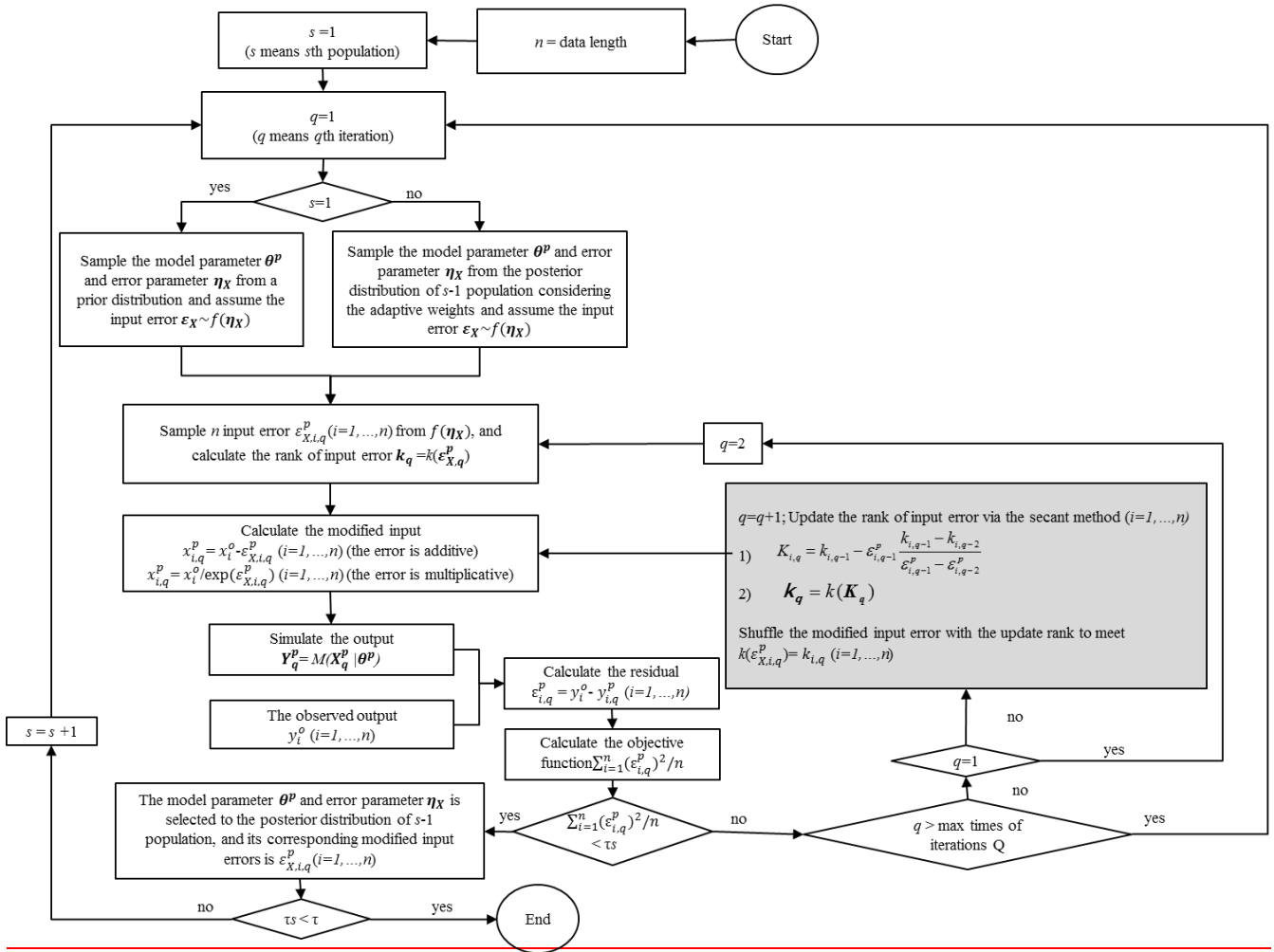
## Acknowledgments

This work was supported by the Australian Research Council [FT120100269] and the Australian Research Council (ARC) Discovery Award [DP170103959] to Dr. Marshall.

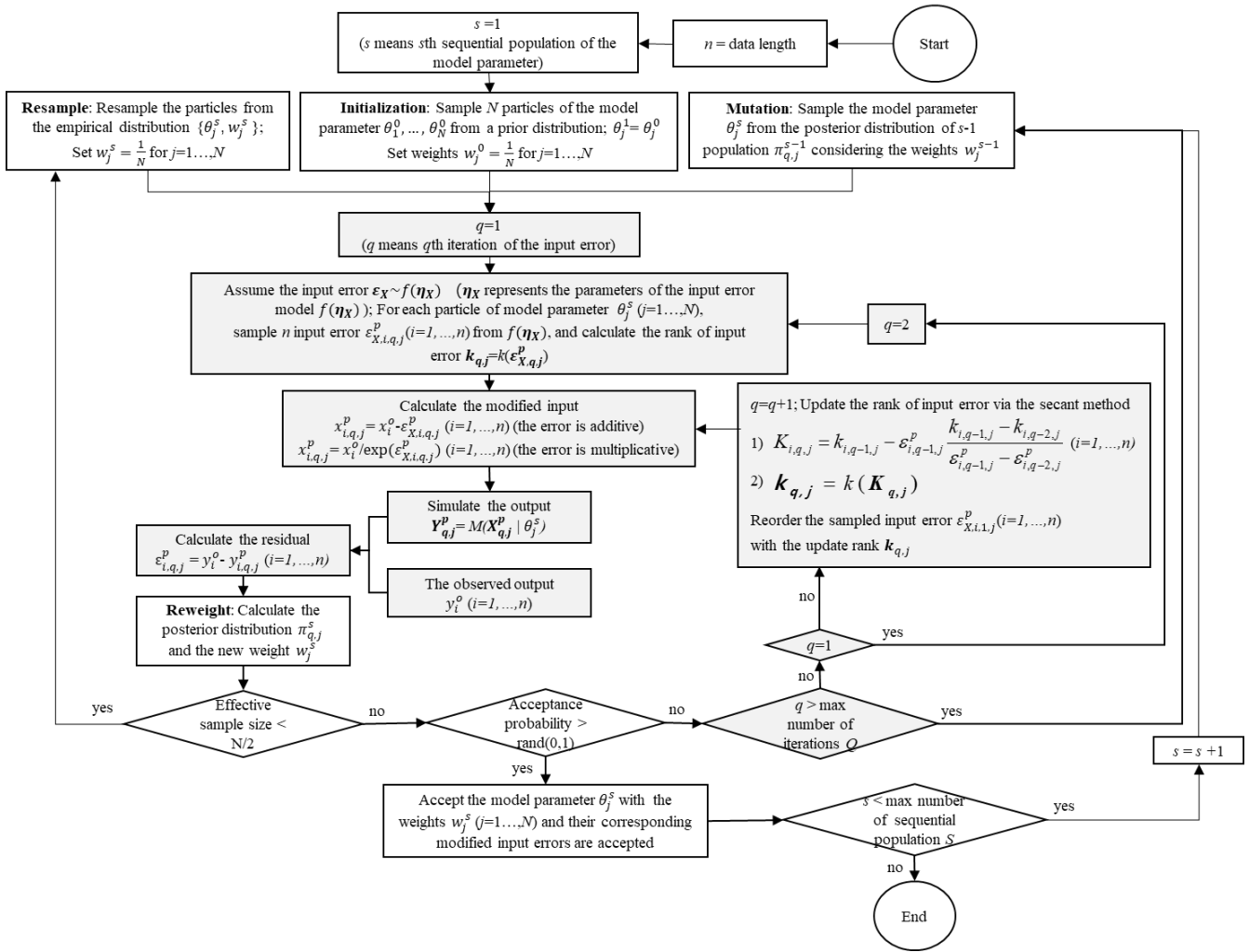
## References

- 705 AJAMI, N. K., DUAN, Q. & SOROOSHIAN, S. 2007. An integrated hydrologic Bayesian multimodel combination framework: Confronting input, parameter, and model structural uncertainty in hydrologic prediction. *Water resources research*, 43.
- BATES, B. C. & CAMPBELL, E. P. 2001. A Markov chain Monte Carlo scheme for parameter estimation and inference in conceptual rainfall - runoff modeling. *Water resources research*, 37, 937-947.
- 710 BEVEN, K. & BINLEY, A. 1992. The future of distributed models: model calibration and uncertainty prediction. *Hydrological processes*, 6, 279-298.
- BONHOMME, C. & PETRUCCI, G. 2017. Should we trust build-up/wash-off water quality models at the scale of urban catchments? *Water research*, 108, 422-431.
- CHAUDHARY, A. & HANTUSH, M. M. 2017. Bayesian Monte Carlo and maximum likelihood approach for uncertainty  
715 estimation and risk management: Application to lake oxygen recovery model. *Water Research*, 108, 301-311.
- DEL MORAL, P., DOUCET, A. & JASRA, A. 2006. Sequential monte carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68, 411-436.
- EVANS, J., WASS, P. & HODGSON, P. 1997. Integrated continuous water quality monitoring for the LOIS river syndromme. *Science of the total environment*, 194, 111-118.
- 720 HAARIO, H., SAKSMAN, E. & TAMMINEN, J. 2005. Componentwise adaptation for high dimensional MCMC. *Computational Statistics*, 20, 265-273.
- HARMEI, R., COOPER, R., SLADE, R., HANEY, R. & ARNOLD, J. 2006. Cumulative uncertainty in measured streamflow and water quality data for small watersheds. *Transactions of the ASABE*, 49, 689-701.

- 725 JEREMIAH, E., SISSON, S., MARSHALL, L., MEHROTRA, R. & SHARMA, A. 2011. Bayesian calibration and uncertainty analysis of hydrological models: A comparison of adaptive Metropolis and sequential Monte Carlo samplers. *Water Resources Research*, 47.
- KAVETSKI, D., KUCZERA, G. & FRANKS, S. W. 2006. Bayesian analysis of input uncertainty in hydrological modeling: 1. Theory. *Water resources research*, 42.
- 730 KLEIDORFER, M., DELETIC, A., FLETCHER, T. & RAUCH, W. 2009. Impact of input data uncertainties on urban stormwater model parameters. *Water Science and Technology*, 60, 1545-1554.
- KUCZERA, G. 1983. Improved parameter inference in catchment models: 1. Evaluating parameter uncertainty. *Water Resources Research*, 19, 1151-1162.
- MARSHALL, L., NOTT, D. & SHARMA, A. 2004. A comparative study of Markov chain Monte Carlo methods for conceptual rainfall-runoff modeling. *Water Resources Research*, 40.
- 735 MCMILLAN, H., KRUEGER, T. & FREER, J. 2012. Benchmarking observational uncertainties for hydrology: rainfall, river discharge and water quality. *Hydrological Processes*, 26, 4078-4111.
- RADWAN, M., WILLEMS, P. & BERLAMONT, J. 2004. Sensitivity and uncertainty analysis for river quality modelling. *Journal of Hydroinformatics*, 6, 83-99.
- 740 RALSTON, M. L. & JENNRICH, R. I. 1978. Dud, A Derivative-Free Algorithm for Nonlinear Least Squares. *Technometrics*, 20, 7-14.
- REFSGAARD, J. C., VAN DER SLUIJS, J. P., HØJBERG, A. L. & VANROLLEGHEM, P. A. 2007. Uncertainty in the environmental modelling process – A framework and guidance. *Environmental Modelling & Software*, 22, 1543-1556.
- 745 RENARD, B., KAVETSKI, D. & KUCZERA, G. 2009. Comment on “An integrated hydrologic Bayesian multimodel combination framework: Confronting input, parameter, and model structural uncertainty in hydrologic prediction” by Newsha K. Ajami et al. *Water Resources Research*, 45.
- RENARD, B., KAVETSKI, D., KUCZERA, G., THYER, M. & FRANKS, S. W. 2010. Understanding predictive uncertainty in hydrologic modeling: The challenge of identifying input and structural errors. *Water Resources Research*, 46.
- 750 RODE, M. & SUHR, U. 2007. Uncertainties in selected river water quality data.
- SARTOR, J. D. & BOYD, G. B. 1972. *Water pollution aspects of street surface contaminants*, US Government Printing Office.
- SCHAEFLI, B., TALAMBA, D. B. & MUSY, A. 2007. Quantifying hydrological modeling errors through a mixture of normal distributions. *Journal of Hydrology*, 332, 303-315.
- 755 SIKORSKA, A. E., DEL GIUDICE, D., BANASIK, K. & RIECKERMANN, J. 2015. The value of streamflow data in improving TSS predictions—Bayesian multi-objective calibration. *Journal of hydrology*, 530, 241-254.
- SMITH, T., SHARMA, A., MARSHALL, L., MEHROTRA, R. & SISSON, S. 2010. Development of a formal likelihood function for improved Bayesian inference of ephemeral catchments. *Water Resources Research*, 46.
- 760 STUBBLEFIELD, A. P., REUTER, J. E., DAHLGREN, R. A. & GOLDMAN, C. R. 2007. Use of turbidometry to characterize suspended sediment and phosphorus fluxes in the Lake Tahoe basin, California, USA. *Hydrological Processes*, 21, 281-291.
- WILLEMS, P. 2008. Quantification and relative comparison of different types of uncertainties in sewer water quality modeling. *Water Research*, 42, 3539-3551.
- 765 WU, X., MARSHALL, L. & SHARMA, A. 2021. Quantifying input error in hydrologic modeling using the Bayesian Error Analysis with Reordering (BEAR) approach. *Journal of Hydrology*, 126202.
- YADAV, M., WAGENER, T. & GUPTA, H. 2007. Regionalization of constraints on expected watershed response behavior for improved predictions in ungauged basins. *Advances in Water Resources*, 30, 1756-1774.

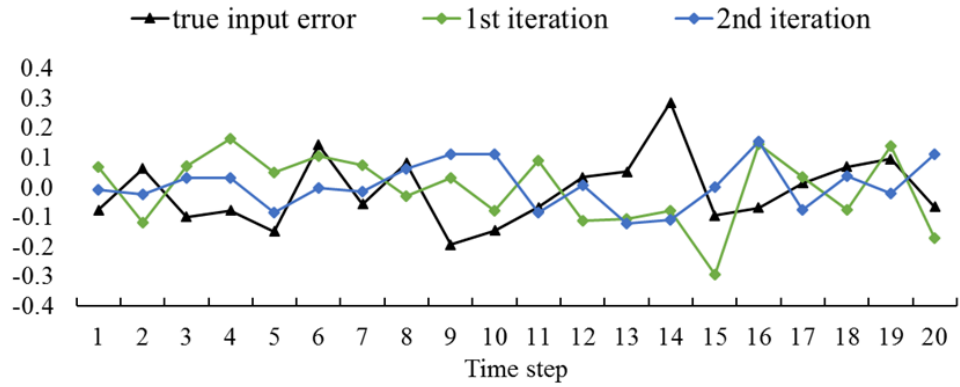


**Figure 1: Flowchart of the algorithm to quantify the input errors via Bayesian error analysis with reshuffling (BEAR) method**

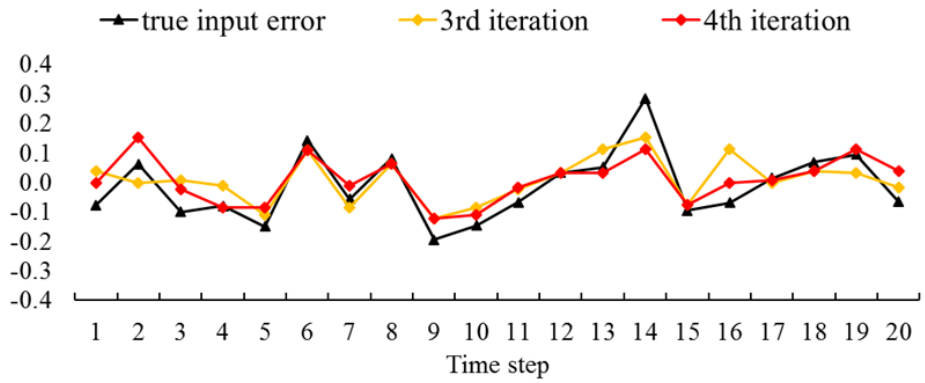


**Figure 1** Flowchart of the algorithm to quantify the input errors via Bayesian error analysis with reordering (BEAR) method in the SMC calibration scheme (The grey charts demonstrate the BEAR method while the white charts demonstrate the SMC algorithm. The details of the BEAR method can refer to Appendix A. The details of the SMC algorithm can refer to the study of Jeremiah et al. (2011), including the Mutation step, the Reweight step and calculating the acceptance probability. rand(0,1) means a number randomly sampled from 0 to 1.)

Randomly  
sampled error

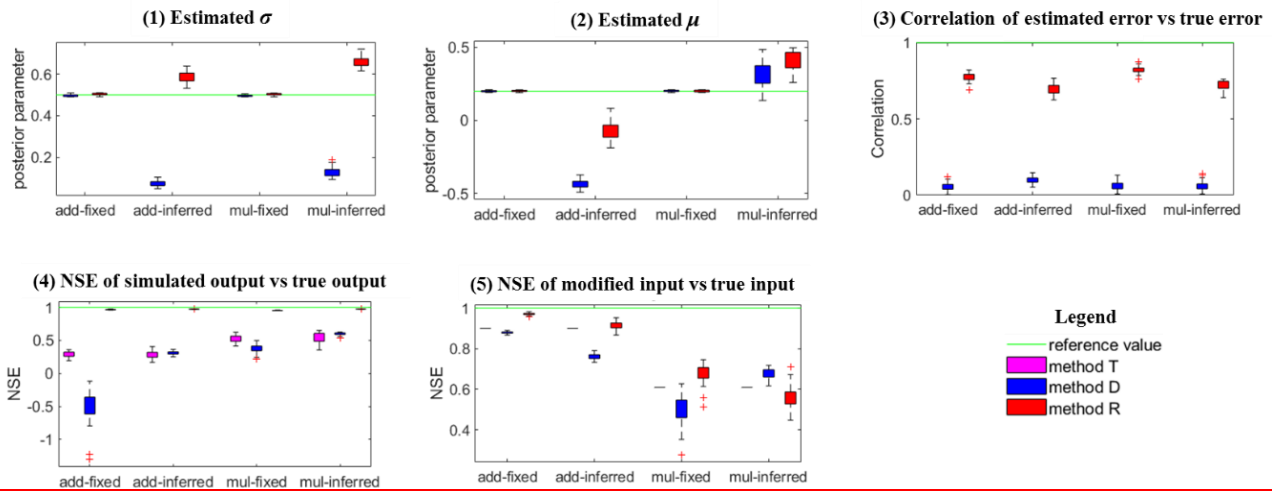


Reshuffled error via  
the secant method



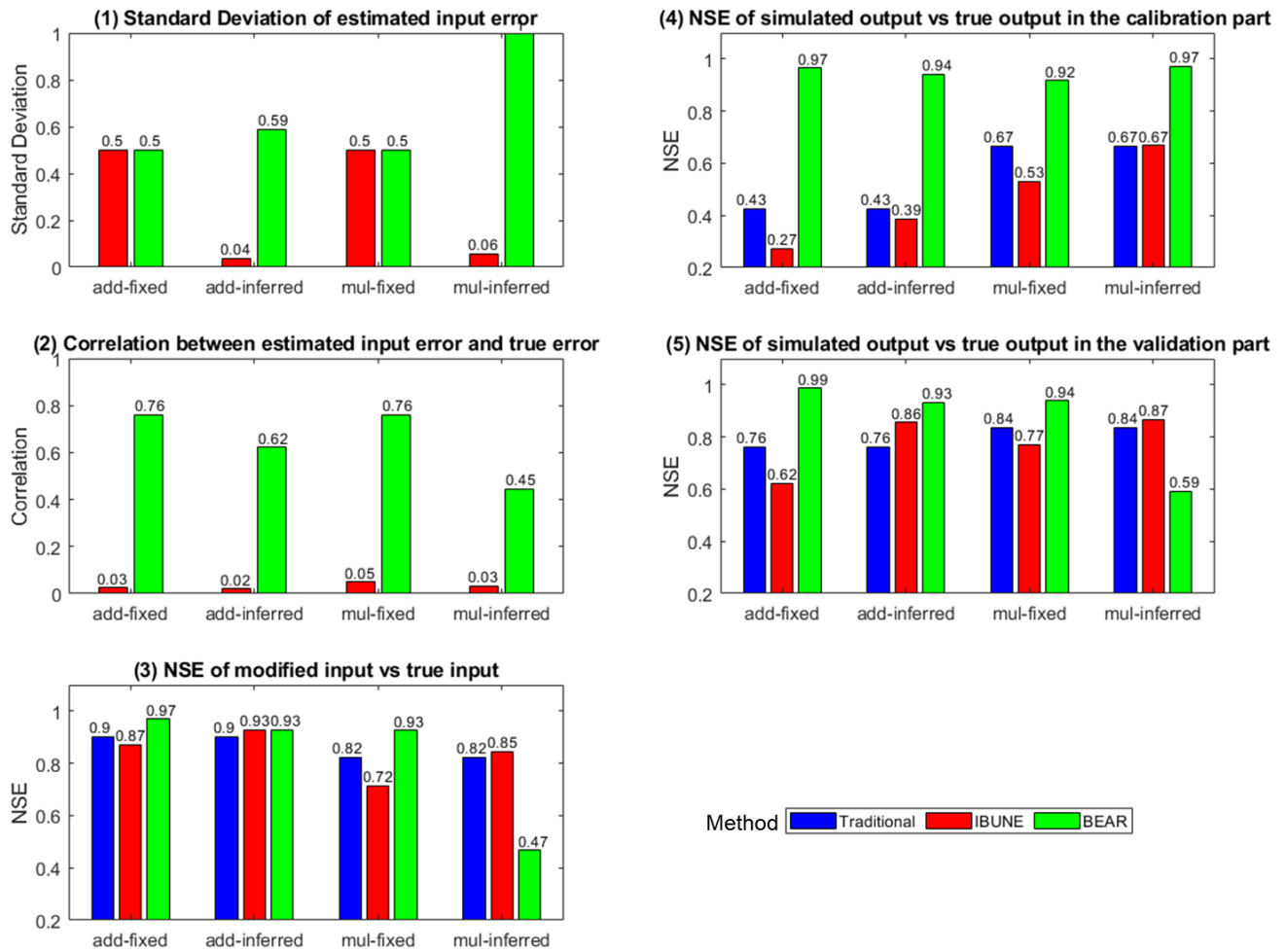
780

**Figure 2: Demonstration of the results in Table 1 before and after reshuffling the errors via the secant method**

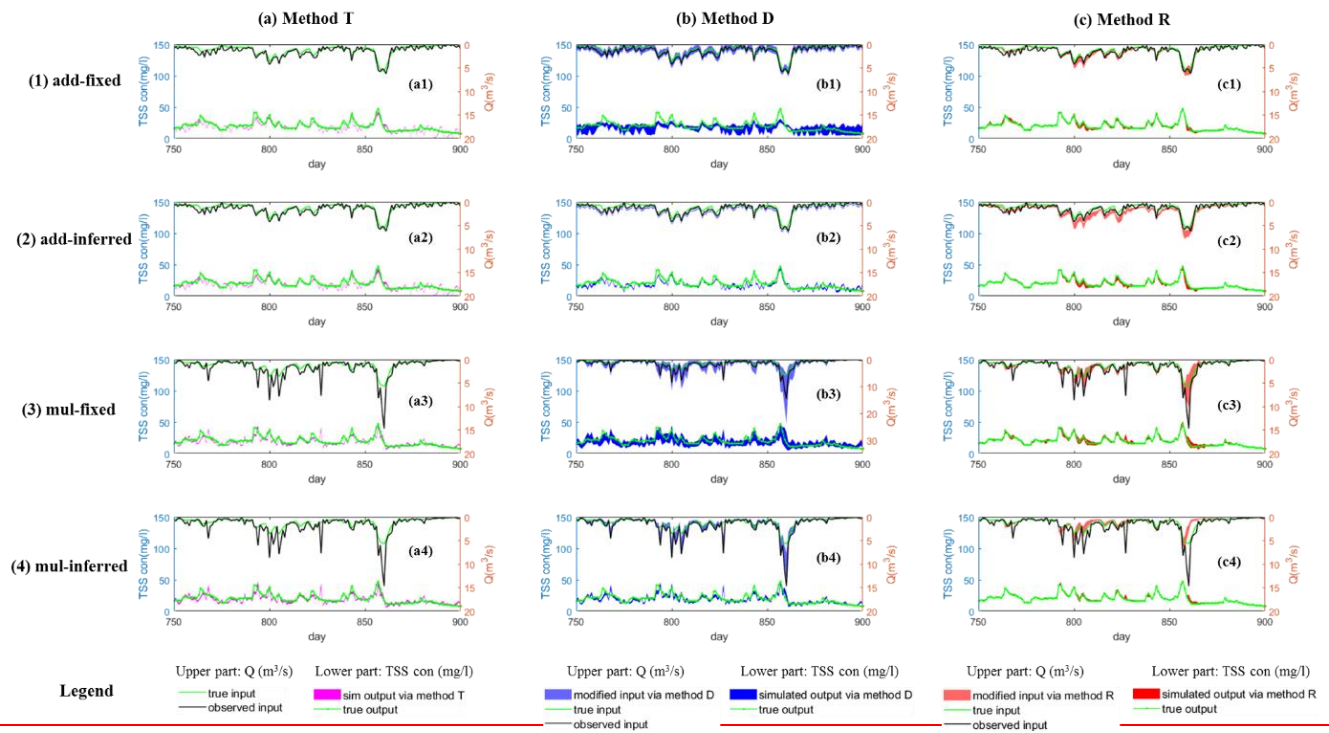


**Figure 3: Comparison of statistical characteristics of four calibration scenarios in the synthetic case (including *add-fixed*, *add-inferred*, *mul-fixed* and *mul-inferred*; notations are given in Table 3) via three calibration methods (including method T, method D and method R, their algorithms are explained in Sect. 2.5)**

785



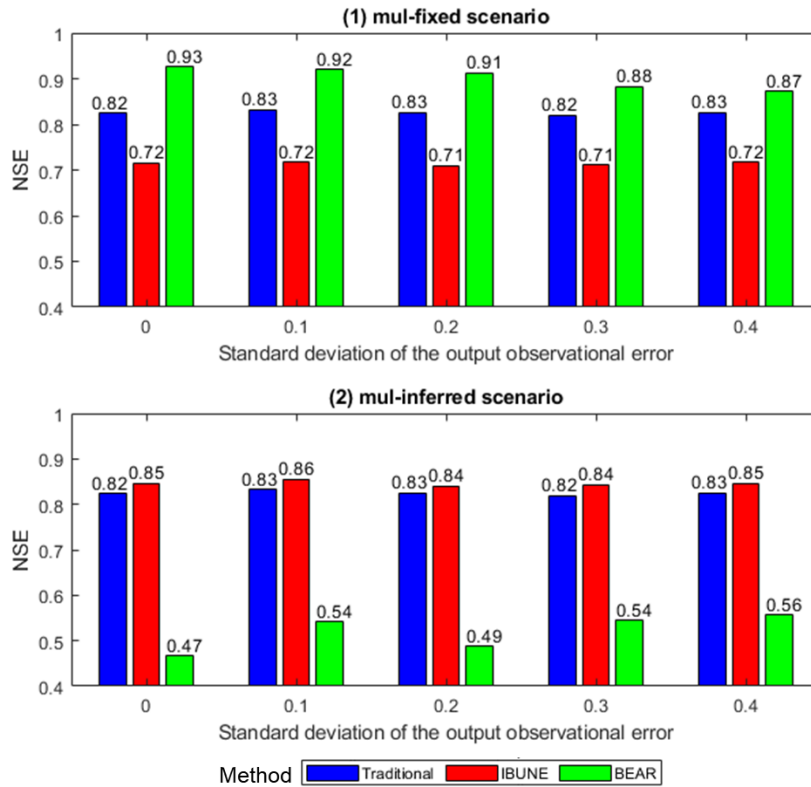
**Figure 2:** Comparison of statistical characteristics of four calibration scenarios in the synthetic case 1 (including *add-fixed*, *add-inferred*, *mul-fixed* and *mul-inferred*; notations are given in Error! Reference source not found.) via three calibration methods (including the traditional method, the IBUNE method and the BEAR method, their algorithms are explained in Sect. 2.4)



**Figure 4: Comparison of time series of synthetic data and uncertainty bands estimated via three calibration methods (including method T, method D and method R; algorithms are explained in Sect. 2.5) for a select period of four calibration scenarios in the synthetic case (including *add-fixed*, *add-inferred*, *mul-fixed* and *mul-inferred*; notations are given in Table 3)**

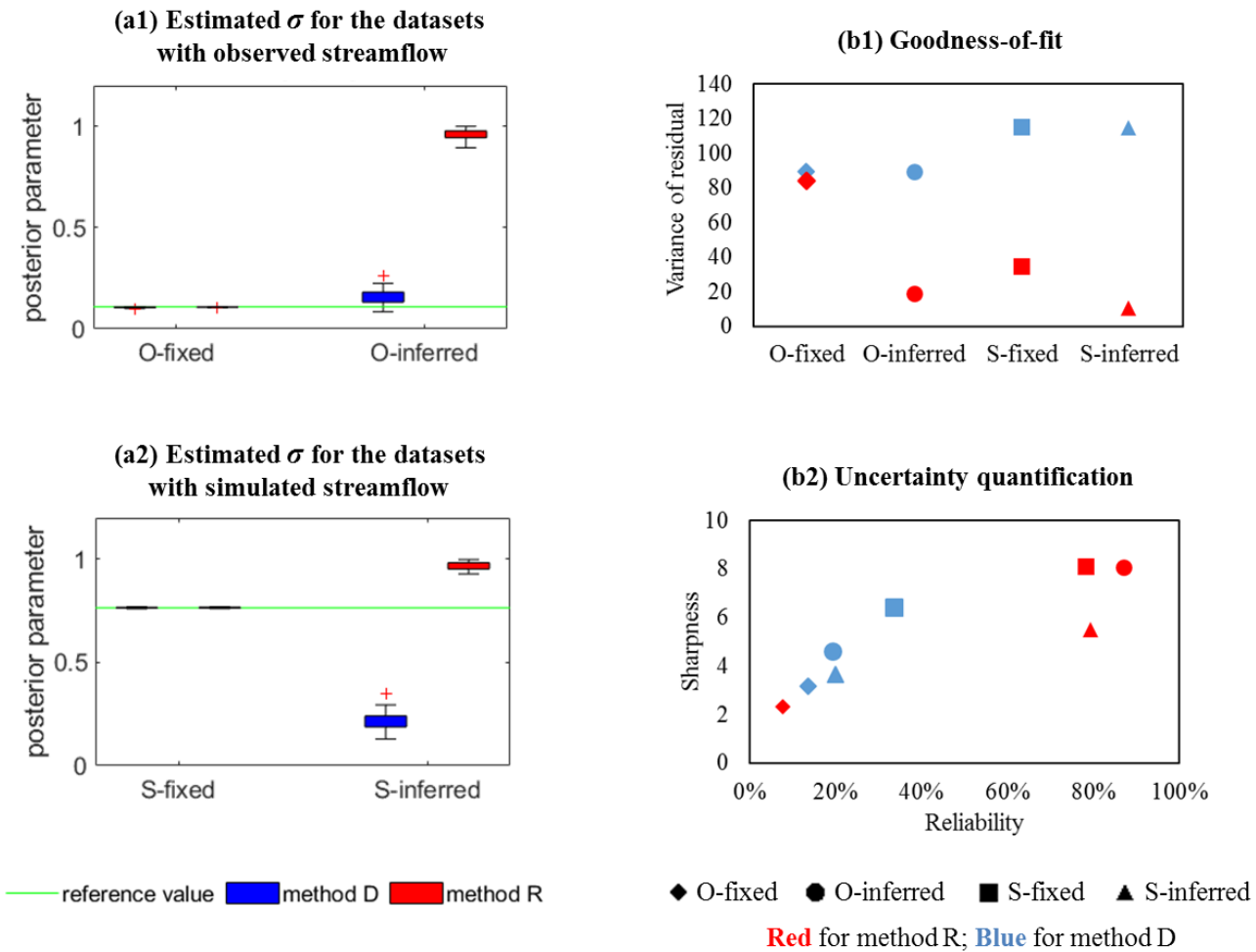
795



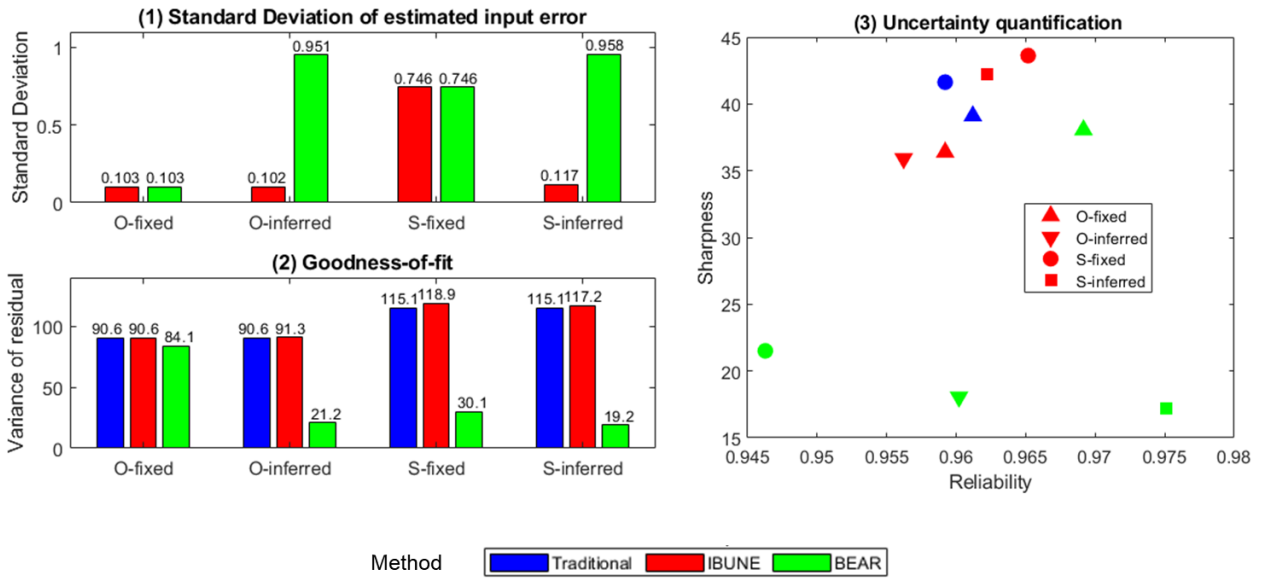


**Figure 3** Comparison of Nash-Sutcliffe efficiency (NSE) of the modified input v.s true input under the interference of the output observational errors with the increasing standard deviations in two calibration scenarios in the synthetic case 2 (including *mul-fixed* and *mul-inferred*; notations are given in Error! Reference source not found.) via three calibration methods (including the traditional method, the IBUNE method and the BEAR method, their algorithms are explained in Sect. 2.4)

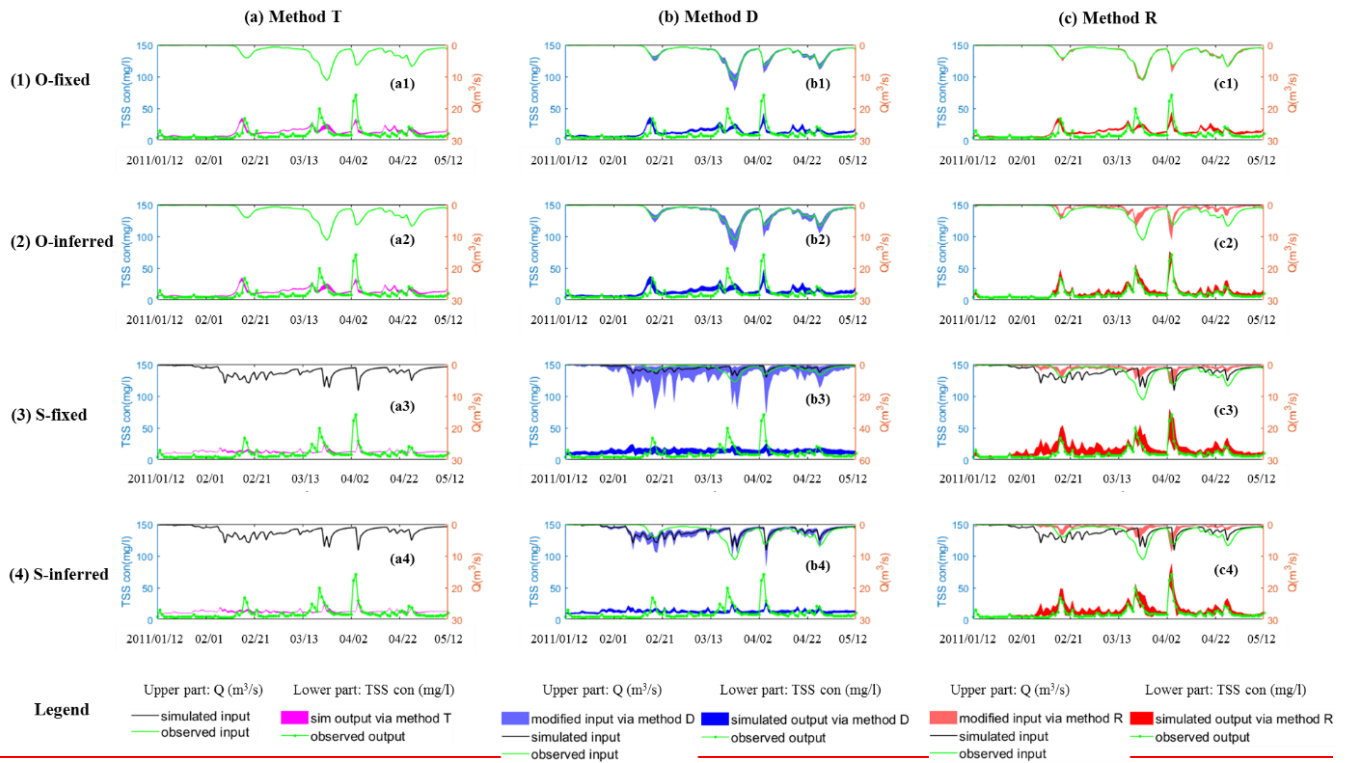
800



805 **Figure 5: Comparison of statistical characteristics of four calibration scenarios in the real case (including *O-fixed*, *O-inferred*, *S-fixed* and *S-inferred*, their notations are given in Table 3) via three calibration methods (including method T, method D and method R, their algorithms are explained in Sect. 2.5)**



810 **Figure 4: Comparison of statistical characteristics of four calibration scenarios in the real case (including *O-fixed*, *O-inferred*, *S-fixed* and *S-inferred*, their notations are given in Error! Reference source not found.) via three calibration methods (including the traditional method, the IBUNE method and the BEAR method, their algorithms are explained in Sect. 2.4)**



815 **Figure 6: Comparison of time series of real data and uncertainty bands estimated via three calibration methods (including method T, method D and method R, algorithms are explained in Sect. 2.5) for a select period of four calibration scenarios in the real case (including O-fixed, O-inferred, S-fixed and S-inferred, notations are given in Table 3)**

**Table 1-1 An example illustrating the rank updating approach via the secant method**

Time Step $i$	Observed data		1 <sup>st</sup> iteration (random sample)					2 <sup>nd</sup> iteration (random sample)				
	$x_i^o$	$y_i^*$	$e_{i,1}^p$	$k_{i,1}$	$x_{i,1}^p$	$y_{i,1}^p$	$e_{i,1}^p$	$e_{i,2}^p$	$k_{i,2}$	$x_{i,2}^p$	$y_{i,2}^p$	$e_{i,2}^p$
1	2.24	24.1	0.07	13	2.18	23.8	0.29	-0.01	9	2.25	24.0	0.13
2	1.87	23.6	-0.12	3	1.99	24.0	-0.49	-0.02	6	1.90	23.8	-0.23
3	1.37	23.1	0.07	14	1.30	22.5	0.58	0.03	14	1.34	22.6	0.43
4	1.02	22.2	0.16	20	0.86	21.2	0.98	0.03	13	0.99	21.7	0.41
5	0.90	22.2	0.05	12	0.85	21.4	0.78	-0.09	3	0.98	22.0	0.21
6	0.99	21.5	0.10	17	0.89	21.8	-0.29	0.00	10	0.99	22.2	-0.70
7	0.76	21.5	0.07	15	0.69	20.8	0.66	-0.02	8	0.78	21.2	0.23
8	0.87	21.4	-0.03	9	0.90	22.0	-0.59	0.06	16	0.81	21.5	-0.09
9	0.60	21.4	0.03	10	0.57	20.1	1.31	0.11	17	0.49	19.5	1.88
10	0.62	21.3	-0.08	7	0.70	21.0	0.31	0.11	18	0.51	19.8	1.52
11	0.70	21.3	0.09	16	0.61	20.4	0.87	-0.09	4	0.78	21.5	-0.20
12	0.85	21.6	-0.11	4	0.97	22.4	-0.76	0.01	12	0.85	21.8	-0.17
13	1.55	24.2	-0.11	5	1.66	24.7	-0.46	-0.12	1	1.67	24.7	-0.53
14	3.20	27.2	-0.08	6	3.28	27.7	-0.54	-0.11	2	3.31	27.8	-0.60
15	1.91	24.6	-0.29	1	2.21	24.9	-0.25	0.00	11	1.91	24.2	0.43
16	1.51	23.6	0.14	19	1.37	22.8	0.80	0.15	20	1.36	22.9	0.72
17	1.26	22.7	0.03	11	1.23	22.7	0.07	-0.08	5	1.34	23.1	-0.36
18	1.09	22.1	-0.08	8	1.16	22.6	-0.56	0.04	15	1.05	22.2	-0.12
19	1.06	22.0	0.14	18	0.92	21.8	0.23	-0.02	7	1.08	22.5	-0.47
20	0.98	22.4	-0.17	2	1.15	22.8	-0.40	0.11	19	0.87	21.6	0.82
Objective function $\frac{1}{n} \sum_{i=1}^n (\sigma_{i,q})^2$							0.40					
								0.47				

**1-2 An example illustrating the rank updating approach via the secant method**

Time Step <i>i</i>	3 <sup>rd</sup> iteration (the secant method)						4 <sup>th</sup> iteration (the secant method)					
	<del><math>K_{i,3}</math></del>	<del><math>k_{i,3}</math></del>	<del><math>\varepsilon_{X,i,3}^p</math></del>	<del><math>x_{i,3}^p</math></del>	<del><math>y_{i,3}^p</math></del>	<del><math>\varepsilon_{i,3}^p</math></del>	<del><math>K_{i,4}</math></del>	<del><math>k_{i,4}</math></del>	<del><math>\varepsilon_{X,i,4}^p</math></del>	<del><math>x_{i,4}^p</math></del>	<del><math>y_{i,4}^p</math></del>	<del><math>\varepsilon_{i,4}^p</math></del>
1	5.63	6	-0.02	2.21	23.9	0.23	13.17	11	0.00	2.24	24.0	0.15
2	8.76	10	0.00	1.88	23.8	-0.20	34.76	20	0.15	1.72	23.3	-0.20
3	14.00	12	0.01	1.36	22.7	0.34	4.65	7	-0.02	1.39	22.9	0.19
4	8.08	9	-0.01	1.03	21.9	0.24	3.26	4	-0.09	1.11	22.2	-0.07
5	-0.33	2	-0.11	1.01	22.1	0.12	0.72	3	-0.09	0.98	22.0	0.24
6	21.84	17	0.11	0.88	21.7	-0.19	19.51	17	0.11	0.88	21.7	-0.18
7	4.28	4	-0.09	0.85	21.6	-0.14	5.54	9	-0.01	0.77	21.2	0.24
8	17.25	16	0.06	0.81	21.5	-0.08	16.00	16	0.06	0.81	21.5	-0.10
9	-6.12	1	-0.12	0.72	21.0	0.40	-3.32	1	-0.12	0.72	21.0	0.39
10	4.18	3	-0.09	0.70	21.0	0.31	-0.87	2	-0.11	0.73	21.1	0.17
11	6.29	7	-0.02	0.72	21.1	0.22	5.44	8	-0.02	0.71	21.0	0.26
12	14.38	14	0.03	0.82	21.6	-0.03	14.36	13	0.03	0.82	21.6	-0.03
13	30.82	19	0.11	1.44	24.0	0.17	14.54	14	0.03	1.52	24.3	-0.07
14	41.98	20	0.15	3.05	27.5	-0.26	33.77	19	0.11	3.09	27.5	-0.30
15	4.71	5	-0.08	1.99	24.6	0.09	3.46	5	-0.08	1.99	24.5	0.12
16	29.64	18	0.11	1.40	23.1	0.55	11.63	10	0.00	1.52	23.4	0.22
17	10.06	11	0.00	1.26	22.8	-0.11	13.56	12	0.01	1.25	22.8	-0.03
18	16.83	15	0.04	1.05	22.2	-0.14	15.00	15	0.04	1.05	22.2	-0.13
19	14.37	13	0.03	1.02	22.2	-0.27	20.79	18	0.11	0.95	21.9	0.08
20	7.60	8	-0.02	1.00	22.2	0.23	3.80	6	0.04	0.94	22.0	0.44

Objective function  $\frac{1}{n} \sum_{i=1}^n (\varepsilon_{i,q})^2$  0.06 0.04

Note:  ~~$x_{i,q}^p = x_{i,q}^o$~~ ,  ~~$\varepsilon_{X,i,q}^p$~~ ,  ~~$y_{i,q}^p = M(x_{i,q}^p | \theta^p)$~~ ,  $M$  is BwMod with the model parameter  $\theta^p$  ( $a=0.04, b=1.6, \kappa = 0.1, S_{max}=70000$ );

~~$\varepsilon_{i,q}^p = y_i^o - y_{i,q}^p$~~ .

In 1<sup>st</sup> and 2<sup>nd</sup> iteration:  ~~$\varepsilon_{X,i,1}^p$~~  and  ~~$\varepsilon_{X,i,2}^p$~~  are randomly sampled from  $N(0,0.01)$ ;  ~~$k_{i,q} = k(\varepsilon_{X,i,q}^p)$~~ .

In 3<sup>rd</sup> and latter iterations:  ~~$K_{i,q} = k_{i,q-1} - \varepsilon_{i,q-1}^p \frac{k_{i,q-1} - k_{i,q-2}}{\varepsilon_{i,q-1}^p - \varepsilon_{i,q-2}^p}$~~ ;  ~~$k_{i,q} = k(K_{i,q})$~~ ;  ~~$\varepsilon_{X,i,q}^p$~~  is  ~~$\varepsilon_{X,j,2}^p$~~  shuffled with  ~~$k_{i,q}$~~  to meet

~~$k_{i,q} = k(\varepsilon_{X,j,2}^p) = k(\varepsilon_{X,i,q}^p)$~~

**Table 2-1** Descriptions of BwMod parameters

Model	Parameter	Description	Unit	Reference value in the synthetic case	Prior range in the case study
BwMod	$a$	wash-off coefficient	-	0.04	(0, 2)
	$b$	wash-off exponent	-	1.6	(0, 3)
	$\kappa$	sediment accumulate rate	-	0.1	(0, 1)
	$S_{max}$	maximum amount of sediment possible to be accumulated	kg	7000	(0, 15000)

825

**Table 3-2** Summary of the calibration scenarios in case studies

Scenario in the synthetic case	Notation	Input error model in the synthetic data generation	Prior information of input error model in calibration
1	<i>add-fixed</i>	<del><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></del>	<del><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></del>
2	<i>add-inferred</i>	<del><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></del>	<del><math>X^o = X^* + \epsilon, \epsilon \sim N(\mu, \sigma^2), \mu \in (-0.5, 0.5), \sigma \in (0, 5)</math></del>
3	<i>mul-fixed</i>	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></del>	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></del>
4	<i>mul-inferred</i>	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></del>	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(\mu, \sigma^2), \mu \in (-0.5, 0.5), \sigma \in (0, 5)</math></del>
1	<i>O-fixed</i>	Observations from the rating curve (USGS database)	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0, \sigma^2), \sigma \in (0.10, 0.11)</math></del>
2	<i>O-inferred</i>		<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0, \sigma^2), \sigma \in (0, 1)</math></del>
3	<i>S-fixed</i>	Simulations from a hydrological model	<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0, \sigma^2), \sigma \in (0.76, 0.77)</math></del>
4	<i>S-inferred</i>		<del><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0, \sigma^2), \sigma \in (0, 1)</math></del>
Scenario in the synthetic case 1	Notation	Input error model in the synthetic input generation	Prior information of input error model in calibration
1	<i>add-fixed</i>	<u><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></u>	<u><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></u>
2	<i>add-inferred</i>	<u><math>X^o = X^* + \epsilon, \epsilon \sim N(0.2, 0.5^2)</math></u>	<u><math>X^o = X^* + \epsilon, \epsilon \sim N(\mu, \sigma^2), \mu=0.2, \sigma \in (0, 1)</math></u>
3	<i>mul-fixed</i>	<u><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></u>	<u><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></u>
4	<i>mul-inferred</i>	<u><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(0.2, 0.5^2)</math></u>	<u><math>X^o = X^* \exp(\epsilon), \epsilon \sim N(\mu, \sigma^2), \mu=0.2, \sigma \in (0, 1)</math></u>
Scenario in the synthetic case 2	Notation	Observational error model in the synthetic output generation	Prior information of input error model in calibration

<u>1</u>	<u><i>mul-fixed</i></u>	<u><math>Y^o = Y^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0, \sigma_Y^2)</math></u>	<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0.2, 0.5^2)</math></u>
<u>2</u>	<u><i>mul-inferred</i></u>	<u><math>\sigma_Y^2 = 0, 0.1, 0.2, 0.3, 0.4</math></u>	<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(\mu, \sigma^2), \mu = 0.2, \sigma \in (0, 1)</math></u>
<u>Scenario in the real case</u>	<u>Notation</u>	<u>Input data source in the real case</u>	<u>Prior information of input error model in calibration</u>
<u>1</u>	<u><i>O-fixed</i></u>	<u>Observations from the rating curve (USGS database)</u>	<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0, \sigma^2), \sigma = 0.103</math></u>
<u>2</u>	<u><i>O-inferred</i></u>		<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0, \sigma^2), \sigma \in (0, 1)</math></u>
<u>3</u>	<u><i>S-fixed</i></u>	<u>Simulations from a hydrological model</u>	<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0, \sigma^2), \sigma = 0.764</math></u>
<u>4</u>	<u><i>S-inferred</i></u>		<u><math>X^o = X^* \exp(\boldsymbol{\varepsilon}), \boldsymbol{\varepsilon} \sim N(0, \sigma^2), \sigma \in (0, 1)</math></u>



**Table 4.3** Characteristics of the study catchments and calibration data

USGS station number	location	State	Drainage area (km <sup>2</sup> )
04087030	Menomonee River at Menomonee Fall	Wisconsin, USA	89.83
land use			Number of Data (days)
Urban (percent)	Agricultural (percent)	Natural (percent)	Period of Data
35	38	27	2009/10/01 - 2012/09/29

830