

1 **Development of Integrated Approaches for Hydrological Data Assimilation**
2 **through Combination of Ensemble Kalman Filter and Particle Filter**
3 **Methods**

4
5 **Y.R. Fan^a, G.H., Huang^{a,b*}, B.W. Baetz^c, Y.P., Li^b, K. Huang^d, X., Chen^e, M. Gao^d**

6
7 ^a Institute for Energy, Environment and Sustainable Communities, University of Regina, Regina,
8 Saskatchewan, Canada S4S 0A2

9 ^b School of Environment, Beijing Normal University, Beijing, China, 100875

10 ^c Department of Civil Engineering, McMaster University, Hamilton, ON L8S 4L8, Canada

11 ^d Faculty of Engineering and Applied Science, University of Regina, Regina, Saskatchewan, Canada S4S
12 0A2

13 ^e State Key Laboratory of Hydrology-Water Resource and Hydraulic Engineering, Hohai University,
14 Nanjing, China, 210098

15

16

17 *Correspondence: Dr. G. H. Huang

18 Institute for Energy, Environment and Sustainable Communities, University of Regina

19 Regina, Saskatchewan, S4S 0A2, Canada

20 Tel: (306) 585-4095

21 Fax: (306) 585-4855

22 E-mail: huangg@uregina.ca

23

24

25

26 **Abstract:**

27 This study improved hydrologic data assimilation through integrating the capabilities of
28 particle filter (PF) and ensemble Kalman filter (EnKF) methods, leading to two integrated
29 data assimilation schemes: the coupled EnKF and PF (CEnPF) and parallelized EnKF and PF
30 (PEnPF) approaches. The applicability and usefulness of CEnPF and PEnPF were
31 demonstrated using a conceptual rainfall-runoff model. The performance of two new
32 developed data assimilation methods and traditional EnKF and PF approaches was tested
33 through a synthetic experiment and two real-world cases with one located in the Jing River
34 basin and one located in the Yangtze river basin. The results show that both PEnPF and
35 CEnPF approaches have more opportunities to provide better results for both deterministic
36 and probabilistic predictions than traditional EnKF and PF approaches. Moreover, the
37 computational time of the two integrated methods is manageable. But the proposed PEnPF
38 may need much more time for some large-scale or time-consuming hydrologic models since
39 it generally needs three times of model runs of EnKF, PF and CEnPF.

40

41 **Keywords:** Hydrologic Prediction, Data assimilation, Ensemble Kalman filter, Particle filter,
42 **Uncertainty**

43

44 **1. Introduction**

45

46 The great increase in computing power and hydrologic data availability has resulted in
47 increasingly use of hydrologic models in real world applications ([Montanari and Brath, 2004](#)).

48 However, significant uncertainties are associated with rainfall-runoff simulation and it is of
49 great importance to account for these uncertainties in hydrologic predictions (e.g.,

50 [Pappenberger and Beven, 2006](#); [Schaake et al., 2006](#); [Brown, 2010](#)). Uncertainty in

51 hydrologic predictions may result from several major sources, including errors in the model
52 structure and model parameters, as well as model initial conditions and forcing data (e.g.,

53 [Ajami et al., 2007](#); [Kavetski et al., 2006a, b](#); [Salamon and Feyen, 2010](#); [Liu et al., 2012](#)).

54 Effective quantification and reduction of these uncertainties is necessary to provide reliable

55 hydrologic forecasts for estimating designated variables in engineering practice, mitigating

56 hydrological risks and improving water resource management policies ([DeChant and](#)

57 [Moradkhani, 2014](#); [Fan et al., 2015a,c](#); [Kong et al., 2015](#); [Li et al., 2015](#); [Yan et al., 2015](#)).

58 Previously, a great number of approaches have been proposed for quantifying the

59 uncertainty in hydrologic predictions ([De Lannoy et al., 2007](#); [Parrish et al., 2012](#); [DeChant](#)
60 [and Moradkhani, 2014](#); [Madadgar and Moradkhani, 2014](#); [Su et al., 2014](#)). Sequential data

61 assimilation techniques are widely used for explicitly dealing with various uncertainties and

62 for optimally merging observations into uncertain model predictions ([Reichle et al., 2002](#);

63 [Moradkhani et al., 2005a](#); [Vrugt et al., 2005](#); [Clark et al., 2008](#); [Xie and Zhang, 2013](#); [Fan et](#)

64 [al., 2015b](#)). The state variables and parameters in a hydrologic model can be continuously

65 updated when new measurements are available through sequential data assimilation

66 techniques, and such a process can highly improve the model predictions. The ensemble

67 Kalman filter (EnKF) and the particle filter (PF) are two of the most widely used sequential

68 data assimilation schemes.

69 The EnKF technique approximates the distribution of the system state using random
70 samples, called ensemble, and replaces the covariance matrix by the sample covariance
71 computed from the ensemble, which is used for state updating in the Kalman filter formula
72 ([Evensen, 1994](#)). The EnKF approach is much attractive in hydrologic predictions due to its
73 features of real-time adjustment and easy implementation ([Reichle et al., 2002](#)). It can
74 provide a general framework for dynamic state, parameter, and joint state-parameter
75 estimation in hydrologic models. For instance, [Moradkhani et al. \(2005a\)](#) initially proposed a
76 dual-state estimation approach based on EnKF for sequential estimation for both the
77 parameters and state variables of a hydrologic model. [Weerts and El Serafy \(2006\)](#) compared
78 the capability of EnKF and particle filter (PF) methods in reducing uncertainty in the
79 rainfall-runoff update and internal model state estimation for flooding forecasting purposes.
80 [Parrish et al. \(2012\)](#) integrated Bayesian model averaging and data assimilation to reduce
81 model uncertainty. [DeChant and Moradkhani \(2014\)](#) combined ensemble data assimilation
82 and sequential Bayesian methods to provide a reliable prediction of seasonal forecast
83 uncertainty. [Shi et al. \(2014\)](#) conducted multiple parameter estimation using multivariate
84 observations via the ensemble Kalman filter (EnKF) for a physically-based land surface
85 hydrologic model. [Pathiraja et al. \(2016a, b\)](#) proposed EnKF-based approaches to detect
86 non-stationary hydrologic model parameters in a paired catchment systems.

87 In comparison with EnKF, the particle filter (PF) method also uses random samples (i.e.
88 particles) to approximate the distributions of the model state. However, these particles are
89 updated forward by using sequential Monte Carlo (SMC) simulation. The most significant
90 advantage of PF is that it relaxes the assumption of Gaussian distribution in state-space model
91 errors, which is required for EnKF. Furthermore, [Liu et al. \(2012\)](#) stated that the PF
92 approaches can reduce numerical instability especially in physically-based or process-based
93 models, since they performs updating on the particle weights instead of the state variables

94 ([Liu et al., 2012](#)). The initial implementation of PF is based on sequential importance
95 sampling, which usually leads to severe deterioration for particles (i.e. only several or even
96 one particle would be available). Consequently, sampling importance resampling (SIR)
97 techniques have been proposed to mitigate this problem (e.g. [Moradkhani et al., 2005b](#); [Li et
98 al., 2015](#); [Fan et al., 2016](#)). However, previous studies in other fields have concluded that the
99 PF method usually requires more samples than other filtering methods and the sample size
100 would increase exponentially with the number of state variables ([Liu and Chen, 1998](#);
101 [Fearnhead and Clifford, 2003](#); [Snyder et al., 2008](#)). Specifically, a great number of samples
102 may be required for reliable characterization of the posterior probability density functions
103 (PDFs) even for small problems with only a few unknown states and parameters ([Liu et al.,
104 2012](#)). Thus, the applications of PF suffer from the number requirement of particles,
105 especially for physically-based distributed hydrologic models ([Liu et al., 2012](#)). Recent
106 improvements for PF are to combine the strengths of sequential Monte Carlo sampling and
107 Markov chain Monte Carlo simulation to achieve a more complete representation of the
108 posterior distribution ([Moradkhani et al., 2012](#); [Vrugt et al., 2013](#)). Such improvements can
109 mitigate sample impoverishment (i.e. a decrease in the diversity of the particles or even a
110 single particle available after resampling steps), and may lead to a more accurate streamflow
111 forecast with small, manageable ensemble sizes ([Moradkhani et al., 2012](#)). Recently, Yan and
112 Moradkhani (2016) demonstrated the application of integration of particle filter and Markov
113 chain Monte Carlo (PF-MCMC) methods by a distributed Sacramento Soil Moisture
114 Accounting (SAC-SMA) model.

115 Both EnKF and PF have been widely used for characterizing uncertainties in hydrologic
116 models. Each of them has its own advantages and drawbacks. The EnKF provides good
117 estimates for very small ensembles but it suffers from its inherent Gaussian assumption ([Shen
118 and Tang, 2015](#)). The PF relaxes the Gaussian assumption and is able to outperform the EnKF

119 if the ensemble size is sufficiently large to prevent filter degeneracy (Moradkhani, 2008;
120 Leisenring and Moradkhani 2012; Shen and Tang, 2015), but it may not recuperate quickly if
121 the particle ensemble consistently over or underestimates the respective observation (Vrugt et
122 al., 2013). Integration of EnKF and PF may be an alternative for overcoming the
123 shortcomings in EnKF and PF, (Frei and Künsch, 2013; Rezaie and Eidsvik, 2012;
124 Plaza-Guingla et al., 2013; Shen and Tang, 2015). For instance, Shen and Tang (2015)
125 proposed a modified ensemble Kalman particle filter for non-Gaussian systems with
126 nonlinear measurement functions by providing a continuous interpolation between the EnKF
127 and PF analysis schemes. The results showed that the proposed method, given an affordable
128 ensemble size, can perform better than the EnKF for nonlinear systems with nonlinear
129 observations (Shen and Tang, 2015).

130 As an extension of previous research, this study aims to develop integrated approaches
131 for hydrologic data assimilation. In detail, two integrated data assimilation approaches are
132 firstly proposed through integrating EnKF and PF: the coupled EnKF and PF (abbreviated as
133 CEnPF) and the parallelized EnKF and PF (abbreviated as PEnPF). The CEnPF sequentially
134 will employ the EnKF and PF to update model parameters and states, in which the EnKF is
135 initially applied to correct model states and parameters, and PF is then adopted to eliminate
136 insignificant particles. In comparison, the PEnPF approach simultaneously updates model
137 states and parameters in parallel through EnKF and PF, and chooses the better estimates as
138 the posterior distributions.

139

140 **2. Methodology**

141 In a sequential data assimilation process, the state variables in a hydrologic model can be
142 evolved forward as follows:

$$143 \quad x_t = f(x_{t-1}, u_{t-1}, \theta) + \omega_{t-1} \quad (1)$$

144 where the subscript t denotes the time step; f is a nonlinear function expressing the system
145 transition from time $t - 1$ to t ; x_t denote the state variables, and θ are the model parameters;
146 ω_{t-1} is considered as process noise (i.e. model error). The model output y_t related to real
147 measurements (e.g. streamflow) can be obtained through the measurement operator $h(\cdot)$,
148 subject to model states and parameters as follows:

$$149 \quad y_t = h(x_t, \theta) + v_t \quad (2)$$

150 where h is the nonlinear function producing forecasted observations; v_t is the observation
151 noise.

152 The essence of the parameter and state estimation problem in the Bayesian filtering
153 framework is to construct the posterior probability density functions (PDFs) of parameters
154 and states conditioned on all previous observations ($y_{1:t-1}$) and current available observation
155 (y_t) (Gordon et al., 1993; Fan et al., 2016). The posterior PDF can be calculated in two steps
156 theoretically: prediction and update, in which the state PDF from the previous state would be
157 integrated through the system model, and the update operation modifies the prediction PDF
158 making use of the latest observations (Han and Li, 2008). The prediction step aims to obtain
159 the prior $p(x_t | y_{1:t-1})$ through the following model:

$$160 \quad p(x_t | y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1} \quad (3)$$

161 where $p(x_t | x_{t-1})$ is the transition probability to describe evolution of states and can be
162 obtained by Equation (1). $p(x_{t-1} | y_{1:t-1})$ is the posterior distribution at time step $t-1$. When
163 new observations at time t are available, the prior can be corrected according to Bayes' rule,
164 formulated as follows:

165
$$p(x_t | y_{1:t}) = \frac{p(y_t | x_t)p(x_t | y_{1:t-1})}{\int p(y_t | x_t)p(x_t | y_{1:t-1})dx_t} \quad (4)$$

166 where $p(x_t | y_{1:t-1})$ represents the prior information; $p(y_t | x_t)$ is the likelihood.

167 The optimal Bayesian solution (i.e. Equations (3) and (4)) is difficult to determine since
168 the evaluation of the integrals may be intractable (Plaza-Guingla et al., 2013). Consequently,
169 approximation methods are applied to address the above issues. Ensemble Kalman filter
170 (EnKF) and PF approaches are the two most widely used methods. The central idea of EnKF
171 and PF is to represent the state probability density function (pdf) as a set of random samples
172 and the difference between these two methods lies in the way of recursively generating an
173 approximation to the state PDF (Weerts and EI Serafy, 2005).

174

175 **2.1. Ensemble Kalman Filter**

176

177 The EnKF and its variants use ensembles of states to approximate the covariance matrices to
178 achieve suboptimal state estimations in which the error statistics are analyzed by numerically
179 solving the Fokker-Planck equation using the Monte Carlo method (Evensen, 2003; Shen and
180 Tang, 2015). EnKF-based filters normally distributed errors and the Monte Carlo approach is
181 applied to approximate the error statistics, as well as compute an approximate Kalman gain
182 matrix for updating model and state variables. A general framework of EnKF for states and
183 parameters updating is described below, followed the description in Moradkhani et al.
184 (2005b).

185

186 In the implementation of EnKF, the prior and posterior distributions for model parameters

187 and state variables are characterized by random samples name “ensembles”. At any given
 188 time t , the prior and posterior distributions of states and parameter are assumed to be denoted
 189 through a set of ensembles below

$$190 \quad X_t^f = (x_{t,1}^f, \dots, x_{t,i}^f, \dots, x_{t,ne}^f) \quad (1)$$

$$191 \quad \Psi_t^f = (\theta_{t,1}^f, \dots, \theta_{t,i}^f, \dots, \theta_{t,ne}^f) \quad (2)$$

$$192 \quad X_t^a = (x_{t,1}^a, \dots, x_{t,i}^a, \dots, x_{t,ne}^a) \quad (3)$$

$$193 \quad \Psi_t^a = (\theta_{t,1}^a, \dots, \theta_{t,i}^a, \dots, \theta_{t,ne}^a) \quad (4)$$

194 where the superscript f indicates the “forecast” values indicating the prior distributional
 195 information and the superscript a indicates the “analyzed” values after assimilation which
 196 denotes the posterior distributional information; the subscript i refers to the i^{th} ensemble
 197 member, and ne denotes the total number of ensembles. Consider a stochastic dynamic-state
 198 model $f(x, u, \theta)$ described by state vector x , parameter vector θ and forcing data u , the state
 199 propagation can be expressed as:

$$200 \quad x_{t+1,i}^f = f(x_{t,i}^a, u_{t,i}, \theta_{t+1,i}^f) + \omega_{t,i}, \quad i = 1, 2, \dots, ne \quad (5)$$

201 where ω_t is the model error term, which follows a Gaussian distribution with zero mean and
 202 covariance matrix P_t . To implement model (5), parameter evolution should be conducted. A
 203 number of parameter evolution approaches have been developed (e.g. [Fan et al., 2015b](#);
 204 [Pathiraja et al., 2016a,b](#)). Among these methods, the random walk method is widely used, in
 205 which stochastic perturbations with mean values of zero and heteroscedastic variances are
 206 added to the analyzed ensembles in the previous stage as follows:

$$207 \quad \theta_{t+1,i}^f = \theta_{t,i}^a + \tau_{t,i}, \quad \tau_{t,i} \sim N(0, \Sigma_t^\theta) \quad (6)$$

208 where Σ_t^θ is the covariance matrix of the analyzed parameter ensembles at time t .

209

210 Based on the forecasts in model states and parameters, the corresponding observation values

211 can be obtained through an observation equation characterized as:

$$212 \quad y_{t+1,i}^f = h(x_{t+1,i}^f, \theta_{t+1,i}^f) + v_{t+1,i}, v_{t+1,i} \sim N(0, \Sigma_{t+1}^y) \quad (7)$$

213 where h represents the operator to transfer the states into the observation space, $v_{t+1,i}$

214 indicates the random perturbation in model prediction, which is drawn from a normal

215 distribution with a mean value of zero and a covariance of Σ_{t+1}^y . When new observations at

216 time step $t+1$ are available, model states and parameters are corrected by assimilating the

217 observation into modelling process, leading to analyzed ensembles indicating the posterior

218 distributions for model states and parameters. Before assimilating observations, stochastic

219 perturbations are usually added to the observations to account for the uncertainty in

220 measurements. In this process, Gaussian noise is generally employed expressed as:

$$221 \quad y_{t+1,i}^o = y_{t+1,i} + \varepsilon_{t+1,i}, \varepsilon_{t+1,i} \sim N(0, \Sigma_{t+1}^{y^o}) \quad (8)$$

222 where $y_{t+1,i}$ represents the raw observation and $\Sigma_{t+1}^{y^o}$ denotes the error covariance. Through

223 assimilating the observations, the posterior states and parameters can be updated by the

224 Kalman update equations:

$$225 \quad x_{t+1,i}^a = x_{t+1,i}^f + K_{xy} [y_{t+1,i}^o - y_{t+1,i}^f] \quad (9)$$

$$226 \quad \theta_{t+1,i}^a = \theta_{t+1,i}^f + K_{\theta y} [y_{t+1,i}^o - y_{t+1,i}^f] \quad (10)$$

227 where K_{xy} , $K_{\theta y}$ are Kalman matrix for states and parameters, which can be expressed as

228 follows (DeChant and Moradkhani, 2012; Pathiraja et al., 2016a):

$$229 \quad K_{xy} = \Sigma_{t+1}^{xy} (\Sigma_{t+1}^y + \Sigma_{t+1}^{y^o})^{-1} \quad (11)$$

$$230 \quad K_{\theta y} = \Sigma_{t+1}^{\theta y} (\Sigma_{t+1}^y + \Sigma_{t+1}^{y^o})^{-1} \quad (12)$$

231 where Σ_{t+1}^{xy} is the cross covariance of the forecasted states $x_{t+1,i}^f$ and the simulated
232 observation $y_{t+1,i}^f$; $\Sigma_{t+1}^{\theta y}$ is the cross covariance between model parameters $\theta_{t+1,i}^f$ and the
233 simulated observation $y_{t+1,i}^f$

234

235 2.2. Particle Filter

236 The PF, similar to the EnKF, is a kind of sequential Monte Carlo method that calculates the
237 posterior distribution of states and parameters by a set of random samples. But PF and its
238 variants are different from EnKF since the ensemble members (or the particles) are not
239 modified, but are combined with different weights (Shen and Tang, 2015). It was found that
240 PF outperforms EnKF by relaxing the assumption of a Gaussian error structure, which allows
241 PF to accurately predict the posterior distribution in the presence of skewed distributions
242 (Moradkhani et al., 2005a; DeChant and Moradkhani, 2012).

243

244 In detail, consider ne independent and identically distributed random variables $x_{t,i} \sim p(x_t | y_{1:t})$
245 for $i = 1, 2, \dots, ne$, the posterior density, based on the sequential importance sampling (SIS)
246 method, can then be approximated as a discrete function:

$$247 \quad p(x_t | y_{1:t}) = \sum_{i=1}^{ne} w_{t,i} \delta(x_t - x_{t,i}) \quad (13)$$

248 where $w_{t,i}$ is the posterior (updated) normalized weight of the i th particle drawn from the
249 proposed distribution; δ is the Dirac delta function. Assume the system state to be a Markov
250 process, and apply the Bayesian recursive expression to the filtering problem. The updating
251 expression for the importance weights (not normalized) is expressed as:

252
$$w_{t,i}^{a*} = w_{t,i}^f \cdot \frac{L_\theta(y_t | x_{t,i}^f) p_\theta(x_{t,i}^f | x_{t-1,i}^f)}{q_\theta(x_{t,i}^f | x_{t-1,i}^f, y_t^f)} \quad (14)$$

253 where $w_{t,i}^f$ is the prior weight, which is equal to the posterior weight at the previous time
 254 step. $w_{t,i}^{a*}$ is the unnormalized posterior weight. Through Equation (14), the importance
 255 weights are sequentially updated when an appropriate proposal distribution $q_\theta(x_{t,i}^f | x_{t-1,i}^f, y_t^f)$ is
 256 given. Consequently, the expression of the proposal distribution will significantly affect the
 257 efficiency and complexity of the PF method. [Gordon et al. \(1993\)](#) have suggested to set

258
$$q_\theta(x_{t,i}^f | x_{t-1,i}^f, y_t^f) = p_\theta(x_{t,i}^f | x_{t-1,i}^f),$$
 resulting in a simplified expression for importance weights:
 259
$$w_{t,i}^a = w_{t,i}^f L_\theta(y_t | x_{t,i}^f) \quad (15)$$

260 Therefore, the normalized updating weight can then be obtained via the following equation:

261
$$w_{t,i}^a = \frac{w_{t,i}^f L_\theta(y_t | x_{t,i}^f)}{\sum_{i=1}^{ne} w_{t,i}^f L_\theta(y_t | x_{t,i}^f)} \quad (16)$$

262 $w_{t,i}^a$ is the normalized posterior weight. $L_\theta(y_t | x_{t,i}^f)$ is the posterior likelihood function. The
 263 choice of an adequate likelihood function has been the subject of considerable debate in
 264 hydrologic and statistics literature ([Vrugt et al., 2013](#)). In the data assimilation process
 265 through PF, the Gaussian likelihood is widely used in a number of fields ([Moradkhani et al.,](#)
 266 [2005b](#); [Weerts and El Serafy, 2006](#); [Salamon and Feyen, 2010](#); [Fan et al., 2016](#)).

267 Consequently, this study will also adopt the Gaussian likelihood expressed as:

268
$$L_\theta(y_t | x_{t,i}^f) = \frac{1}{\sqrt{2\pi R_t}} \exp\left(-\frac{1}{2R_t} [y_t - y_{t,i}^f]^2\right) \quad (17)$$

269

270 For the particle filter through SIS, a serious limitation is the depletion of the particle set,
 271 which means that, after a few iterations (time steps), all the particles except one are discarded

272 because their importance weights are insignificant (Doucet, et al. 2001). To address the above
273 issue, sampling importance resampling (SIR) algorithms are usually applied to eliminate the
274 particles with small importance weights and replace them by particles with large importance
275 weights. A number of resampling approaches have been developed, such as multinomial
276 resampling, systematic resampling, residual resampling, and grouping-based resampling
277 approaches (Li et al., 2015)

278

279 **2.3. Integration of EnKF and PF for Hydrologic Data Assimilation**

280

281 The application of EnKF is constrained by its assumption of Gaussian errors while the PF
282 requires a large sample size for providing reliable predictions. In this study, we extend the
283 previous research to provide two integrated data assimilation schemes: the coupled EnKF and
284 PF (abbreviated as CEnPF) and the parallelized EnKF and PF (abbreviated as PEnPF)
285 approaches to characterize uncertainty in hydrologic models.

286

287 **2.3.1. the coupled EnKF and PF (CEnPF) approach**

288 The CEnPF sequentially uses the EnKF and PF to update model parameters and states, in
289 which EnKF is first applied to correct model states and parameters, and PF is then adopted to
290 eliminate insignificant particles (see Figure 1). The detailed procedures for the
291 implementation of CEnPF are presented as follows:

292 *Step 1.* Similar to the implementation of EnKF and PF, the model initial conditions should be
293 assumed before implementing CEnPF. In this study, the initial state variables and parameters

294 are sampled from the corresponding uniform distributions:

$$295 \quad x_{1,i} \sim U(x^L, x^U), i = 1, 2, \dots, ne, \quad x \in R^{N_x} \quad (18)$$

$$296 \quad \theta_{1,i} \sim U(\theta^L, \theta^U), i = 1, 2, \dots, ne, \quad \theta \in R^{N_\theta} \quad (19)$$

297 *Step 1.* Assign prior weights for the ensembles. In general, the prior weights are assigned
 298 uniformly as follows:

$$299 \quad w_{t,i} = 1/ne, i = 1, 2, \dots, ne \quad (20)$$

300 *Step 3.* At any time step t , model states at current step can be forecasted based the posterior
 301 states in step $t-1$ and the prior parameters in the current step by using model operator f :

$$302 \quad x_{t,i}^f = f(x_{t-1,i}^a, u_{t,i}, \theta_{t,i}^f) + \omega_{t,i}, \omega_t \sim N(0, \Sigma_t^m), i = 1, 2, \dots, ne \quad (21)$$

303 where parameters $\theta_{t,i}^f$ are obtained by Equation (6).

304 *Step 5.* Observation simulation: Use the observation operator h to propagate the model state
 305 forecast:

$$306 \quad y_{t,i}^f = h(x_{t,i}^f, \theta_{t,i}^f) + v_{t,i}, v_{t+1,i} \sim N(0, \Sigma_t^y), i = 1, 2, \dots, ne \quad (22)$$

307 *Step 6.* Parameters and states updating: Update the parameters and states via the EnKF
 308 updating equations

$$309 \quad x_{t,i}^a = x_{t,i}^f + K_{xy}[y_{t,i}^o - y_{t,i}^f] \quad (23)$$

$$310 \quad \theta_{t,i}^a = \theta_{t,i}^f + K_{\theta y}[y_{t,i}^o - y_{t,i}^f] \quad (24)$$

311 where $x_{t,i}^a$ and $\theta_{t,i}^a$ are the updated state and parameter values and K_{xy} and $K_{\theta y}$ are the
 312 Kalman matrix for states and parameters obtained by Equations (11) and (12).

313 *Step 7.* Estimate the likelihood:

$$314 \quad L(y_t | x_{t,i}^a, \theta_{t,i}^a) = \frac{1}{\sqrt{2\pi R_t}} \exp\left(-\frac{1}{2R_t} [y_{t,i}^o - h(x_{t,i}^a, \theta_{t,i}^a)]^2\right) \quad (25)$$

$$315 \quad p(y_t | x_{t,i}^a, \theta_{t,i}^a) = \frac{L(y_t | x_{t,i}^a, \theta_{t,i}^a)}{\sum_{i=1}^{ne} L(y_t | x_{t,i}^a, \theta_{t,i}^a)} = p(y_{t,i}^o - h(x_{t,i}^a, \theta_{t,i}^a) | R_t) \quad (26)$$

316 *Step 8.* update weight for the analyzed ensemble values:

$$317 \quad w_{t,i}^a = \frac{w_{t,i}^f \cdot p(y_{t,i}^o - h(x_{t,i}^a, \theta_{t,i}^a) | R_t)}{\sum_{i=1}^{ne} w_{t,i}^f \cdot p(y_{t,i}^o - h(x_{t,i}^a, \theta_{t,i}^a) | R_t)} \quad (27)$$

318 where $w_{t,i}^f$ are the prior sample weights and are usually set to be $1/ne$.

319 *Step 9.* Resampling: Apply resampling procedure proposed by [Moradkhani et al. \(2005a\)](#) to
320 eliminate the abnormal samples in $x_{t,i}^a$, and $\theta_{t,i}^a$, and generate resampled ensembles denoted

321 as $x_{t-resamp,i}^a, \theta_{t-resamp,i}^a$.

322 *Step 10.* Parameter perturbation: take parameter evolution to the next stage through adding
323 small stochastic error around the sample:

$$324 \quad \theta_{t+1,i}^f = \theta_{t-resamp,i}^a + \varepsilon_{t,i}, \quad \varepsilon_{t,i} \sim N(0, \eta S(\theta_{t-resamp,i}^a)) \quad (28)$$

325 where η is a hyper-parameter which determines the radius around each sample being explored;

326 $S(\theta_{t-resamp,i}^a)$ is the standard deviation of the analyzed ensemble values.

327 *Step 11.* Check the stopping criterion: if measurement data is still available in the next stage, t
328 = $t + 1$ return to step 3; otherwise, stop.

329

330 In CEnPF, model parameters and states are initially updated through Kalman update

331 equations, then the updated states and parameters are corrected again through PF procedure to

332 eliminate abnormal or insignificant state and parameters and replace them by significant ones

333 by sampling importance resampling procedure. Compared with EnKF, the CEnPF can be

334 applicable for nonlinear and non-Gaussian systems. At any time step t , even though the EnKF

335 procedure may not produce optimal states and parameters under nonlinear and non-Gaussian

336 systems, the following PF procedure can remove non-optimal ensembles (i.e. insignificant

337 samples) and replace them with significant ones. In comparison with PF, the proposed CEnPF

338 firstly reduces the sample requirement for large-scale models since the inherent EnKF

339 procedure can achieve satisfactory performance with a moderate sample size; it can also
340 adjust the ensemble values to fit the observations well especially when the particle ensembles
341 consistently over or underestimates the respective observations.

342 **2.3.2. the parallelized EnKF and PF (PEnPF) approach**

343 In comparison with CEnPF, the PEnPF approach simultaneously updates model states and
344 parameters in parallel through EnKF and PF, and chooses the better estimates as the posterior
345 distributions (see Figure 2). The full description of the PEnPF procedures is illustrated as
346 follows:

347 *Step 1.* Model state initialization: Initialize N_x -dimensional model state variables and
348 N_θ -dimensional model parameters from uniform distributions expressed as Equations (18)
349 and (19)

350 *Step 2.* Sample weight assignment: Assign the prior weights uniformly to the particles
351 expressed as Equation (20):

352 *Step 4.* Model state forecast step: Propagate the ne state variables and model parameters
353 forward in time using model operator f by Equation (21).

354 *Step 5.* Observation simulation: Use the observation operator h to propagate the model state
355 forecasts by Equation (22):

356 *Step 6.* Parameters and states updating based on EnKF: This step is further divided into two
357 procedures: model parameters and states are updated by Kalman updating scheme and the
358 updated ensembles are evaluated by a mismatch index proposed by [Gu and Oliver \(2007\)](#).

359 *6a.* Obtain the analyzed estimations through Kalman updating scheme expressed as Equations
360 (23) and (24)

361 *6b.* Evaluate the data match term for the analyzed estimation by the mismatch index
362 expressed by:

$$363 \quad S(x_{t,i}^a, \theta_{t,i}^a) = \sum_{i=1}^{ne} (h(x_{t,i}^a, \theta_{t,i}^a) - y_{t,i}^o)^T R_t^{-1} (h(x_{t,i}^a, \theta_{t,i}^a) - y_{t,i}^o) \quad (29)$$

364 Such an index has been adopted in several data assimilation literatures (e.g. [Gu and Oliver](#)

365 2007; Chen and Oliver, 2013; Zhang et al., 2014) to evaluate history-matching results. In this
 366 study, this index is used to evaluate the performance of the updated states and parameters
 367 obtained from Kalman updating scheme.

368 *Step 7.* Different from the CEnPF in which PF updates model parameters and states based on
 369 the analyzed state and parameter values from EnKF, the PF procedure in PEnPF also update
 370 model states and parameters from the priori states and parameters at time t . Therefore, the
 371 likelihood function can be expressed as:

$$372 \quad L(y_t | x_{t,i}^f, \theta_{t,i}^f) = \frac{1}{\sqrt{2\pi R_t}} \exp\left(-\frac{1}{2R_t} [y_{t,i}^o - h(x_{t,i}^f, \theta_{t,i}^f)]^2\right) \quad (30)$$

$$373 \quad p(y_t | x_{t,i}^f, \theta_{t,i}^f) = \frac{L(y_t | x_{t,i}^f, \theta_{t,i}^f)}{\sum_{i=1}^{ne} L(y_t | x_{t,i}^f, \theta_{t,i}^f)} = p(y_{t,i}^o - h(x_{t,i}^f, \theta_{t,i}^f) | R_t) \quad (31)$$

374 Then, the updated weights denoted as $w_{t,i}^a$ for each particle can be obtained as:

$$375 \quad w_{t,i}^a = \frac{w_{t,i}^f \cdot p(y_{t,i}^o - h(x_{t,i}^f, \theta_{t,i}^f) | R_t)}{\sum_{i=1}^{ne} w_{t,i}^f \cdot p(y_{t,i}^o - h(x_{t,i}^f, \theta_{t,i}^f) | R_t)} \quad (32)$$

376 Based on the updated weights, those particles can be resampled to remove those samples with
 377 insignificant weights. A number of resample methods have been developed and the
 378 multinomial resampling method proposed by Moradkhani et al. (2005a) is used. Therefore,
 379 the resampled particles denoted as $\theta_{t-resamp,i}$ and $x_{t-resamp,i}$ can be obtained. The performance
 380 of the resampled particles is also evaluated by the mismatch index expressed as:

$$381 \quad S(x_{t-resamp,i}, \theta_{t-resamp,i}) = \sum_{i=1}^{ne} (h(x_{t-resamp,i}, \theta_{t-resamp,i}) - y_{t,i}^o)^T R_t^{-1} (h(x_{t-resamp,i}, \theta_{t-resamp,i}) - y_{t,i}^o) \quad (33)$$

382 *Step 8.* Choose the posterior estimations for states and parameters by the following criteria:

383 If $S(x_{t+1-resamp,i}, \theta_{t+1-resamp,i}) \leq S(x_{t+1,i}^a, \theta_{t+1,i}^a)$, $\theta_{t-resamp,i}, x_{t-resamp,i}$ would be the posterior
 384 estimations at current stage; otherwise, $x_{t+1,i}^a$, and $\theta_{t+1,i}^a$ would be the posterior estimations.

385 *Step 9* Parameter perturbation: take parameter evolution to the next stage through add small
 386 stochastic error around the sample (take the EnKF estimation as an example):

$$387 \quad \theta_{t+1,i}^f = \theta_{t,i}^a + \varepsilon_{t,i}, \quad \varepsilon_{t,i} \sim N(0, \eta S(\theta_{t,i}^a)) \quad (34)$$

388 where η is a hyper-parameter which determines the radius around each sample being explored;
389 $S(\theta_{t,i}^a)$ is the standard deviation of the analyzed ensemble values.

390 *Step 10.* Check the stopping criterion: if measurement data is still available in the next stage, t
391 $= t + 1$ return to step 3; otherwise, stop.

392

393 Through PEnPF, the better estimations from EnKF and PF will be chosen as the posterior
394 states and parameters, which may lead to improved predications for model states and
395 simulated observations. Similar to CEnPF, the PEnPF can be applicable for nonlinear and
396 non-Gaussian systems where once the estimates from EnKF are non-optimal, the estimates
397 from PF will be adopted. Also, the ensembles will be adjusted through EnKF when the
398 resulting predictions are consistently over or underestimates the respective observations.

399

400 **3. Synthetic Experiments**

401 **3.1. Rainfall-Runoff Model**

402

403 In this study, the Hymod, is adopted to test the efficiency of the CEnPF and PEnPF
404 approaches. Hymod is a non-linear rainfall-runoff conceptual model which can be run in a
405 minute/hour/daily time step (Moore, 1985). In Hymod, the soil moisture storage is
406 characterized by a spatial probability distribution function and the runoff is routed to the
407 catchment outlet by a fast linear-routing process (nominally event runoff) and a slow
408 nonlinear routing process (nominally baseflow), as shown in Figure 3 (Moore, 2007). A
409 cumulative distribution function (CDF) is proposed to describe such variability of soil
410 moisture capacities, expressed as (Moore, 1985, 2007):

$$411 \quad F(c) = 1 - \left[1 - \frac{c}{C_{max}} \right]^{b_{exp}}, \quad 0 \leq c \leq C_{max} \quad (35)$$

412 where C_{max} [L] is the maximum soil moisture capacity within the catchment and b_{exp} [-] is the
 413 degree of spatial variability of soil moisture capacities and affects the shape of the CDF. Five
 414 parameters are involved in Hymod for calibration based on observations: (i) the maximum
 415 storage capacity (C_{max}), (ii) spatial variability of soil moisture capacity (b_{exp}), (iii) the
 416 partitioning factor between the two series of reservoir tanks (α), (iv) the residence for the
 417 time quick-flow tank (R_q), and (v) the residence time for the slow-tank (R_s). Two inputs are
 418 required to force this model: precipitation, P (mm/day), and potential evapotranspiration, ET
 419 (mm/day).

420

421 -----

422 Place Figure 3 Here

423 -----

424

425 **3.2. Synthetic Experiments**

426

427 In this study, synthetic experiments are initially applied to test the applicability of the CEnPF
 428 and PEnPF approaches. The “true” observations are first defined when the model is run for a
 429 set of meteorological and initial conditions in the synthetic experiment ([Moradkhani, 2008](#)).

430 The “true” model parameters are predefined before the synthetic experiment. The model
 431 inputs, including the potential evapotranspiration, ET (mm/day), and mean areal precipitation,
 432 P (mm/day), are generated based on onsite meteorological data, in which the mean areal
 433 precipitation data are generated based on the rain station measurements in the watershed, and
 434 the potential evapotranspiration values are interpolated based on data from national weather

435 stations near the watershed.

436

437 Stochastic perturbations are required in a data assimilation framework to account for the
438 uncertainties in model inputs, parameters and structures. In the synthetic experiments,
439 random perturbations are added to precipitation and potential evapotranspiration (ET)
440 observations to account for their uncertainties. For potential evapotranspiration, a Gaussian
441 noise distribution is recommended by a number of researchers (e.g. [DeChant and Moradkhani,](#)
442 [2012; Moradkhani et al., 2012; Chen et al., 2013; Rasmussen et al., 2015](#)). For precipitation,
443 some studies have applied Gaussian noise (e.g. [Rasmussen et al., 2015](#)), while other studies
444 have concluded that log-normal noise may perform better (e.g. [DeChant and Moradkhani,](#)
445 [2012; Moradkhani et al., 2012](#)). In this study, the log-normal noise is adopted for the
446 synthetic experiments, while Gaussian noises are employed for potential evapotranspiration,
447 synthetic observations and model predictions. The proportionality factors are set to be 0.2 for
448 all data in the synthetic experiments.

449

450 **3.3. Evaluation Criteria**

451 The root-mean-square error (RMSE), and the Nash-Sutcliffe efficiency (NSE)
452 coefficient will be adopted to evaluate the performance of different data assimilation methods.
453 These two indices also served as the responses in the multi-level factorial design to
454 visualizing the effects of stochastic perturbations. The formations of RMSE and NSE are
455 expressed as follows:

$$456 \quad RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (Q_i - P_i)^2} \quad (36)$$

$$NSE = 1 - \frac{\sum_{i=1}^N (Q_i - P_i)^2}{\sum_{i=1}^N (Q_i - \bar{Q})^2} \quad (37)$$

where N is the total number of observations (or predictions), Q_i are the observed values, P_i are the estimated values, and \bar{Q} is the mean of all observed and estimated values.

460

Both RMSE and NSE merely measure the accuracy of the expected value and show the ability of each data assimilation technique to track the observations (Dechant et al., 2012).

However, they are unable to evaluate the performance of predictive distribution from ensemble forecasts (Renard et al., 2010). Consequently, probabilistic measures are required to further provide a description of ensemble forecasts for different data assimilation schemes. In this study, the continuous ranked probability score (CRPS) and resolution (π) are used, which are formulated as follows (Murphy and Winkler, 1987; Hersbach, 2000; Madadgar et al., 2014):

$$CRPS = \int_{-\infty}^{+\infty} [F^f(x) - F^o(x)]^2 dx \quad (38)$$

where where F^f and F^o are CDFs for forecasts and observations, respectively

$$\pi = \frac{1}{T} \sum_{t=1}^T \frac{E[y_{t,i}]}{\sigma[y_{t,i}]} \quad (39)$$

where $E[y_{t,i}]$ is the expected value of ensemble predictions at time t and $\sigma[y_{t,i}]$ is the standard deviation of ensemble predictions at time t .

474

The CRPS is a measurement of error for probabilistic prediction. A small CRPS value indicates a better model performance, with the value of zero suggesting a perfect accuracy for model prediction. The index of resolution provides a description of precision of ensemble predictions with greater values suggesting larger uncertainty of forecasts (Madadgar et al.,

479 2014)

480

481 3.4. Results Analysis

482

483 To demonstrate the capability of the proposed CEnPF and PEnPF approaches in parameters
484 and state quantification for hydrologic models, synthetic experiments were performed with
485 Hymod. Table 1 shows the “true” parameter set for the synthetic experiments. The initial
486 ensembles for the five parameters (i.e. C_{max} , b_{exp} , α , R_s , R_q) are sampled uniformly from
487 predefined intervals as shown in Table 1. The initial ensembles for the state variable of
488 storage are sampled from a normal distribution with a mean value of 0.05 and a standard
489 deviation to be proportional to the mean value (the proportional factor is set to be 0.1). The
490 initial samples for the slow flow tank are also sampled from a similar normal distribution
491 with a mean value of 2.14. The initial samples for the three quick flow tanks are set to be 0,
492 and the sample size used in the synthetic experiment was 200.

493

494 Figure 4 shows the comparison between the ensemble predictions and the synthetically
495 generated true discharge values obtained from the EnKF, PF, CEnPF and PEnPF approaches.
496 The results indicate that the ensemble means of streamflow predictions from the four
497 methods can track well the observed discharge data. The ranges formulated by 5% and 95%
498 percentiles (i.e. 90% confidence intervals) of streamflow predictions can adequately bracket
499 the observations. In addition, ensemble predictions for two state variables, namely the storage
500 and the flow in the slow tank of Hymod are plotted and compared with their true values in the
501 experiment, as shown in Figure 4. The results show that, for all the four data assimilation
502 schemes, the deterministic predictions (i.e. predictive means in this study) of state variables
503 can well trace the fluctuation of their true values. Moreover, almost all the true values for the

504 two state variables are located in the predictive intervals of the ensemble predictions of the
505 four approaches.

506

507 [Figure 5](#) describes the comparison of the convergence of each parameter from the EnKF, PF,
508 CEnPF and PEnPF approaches. It is observed that identifiability of one parameter depends on
509 the filtering approaches. For instance, all five parameters in Hymod are identifiable if the PF
510 is employed, while in comparison the parameters of C_{max} and b_{exp} are unidentifiable for EnKF.
511 For the two developed methods, CEnPF and PEnPF, the five parameters of Hymod can be
512 well identified by CEnPF. Moreover, compared with PF, the proposed CEnPF can still
513 rejuvenate ensembles in larger spaces than PF, which may lead to more reliable estimations
514 for parameter posterior distributions. In comparison, parameter evolution patterns generated
515 by PEnPF are similar with those from EnKF, which means that C_{max} and b_{exp} are
516 unidentifiable in this data assimilation scheme. This is due to the mechanism of ensembles
517 rejuvenation in PEnPF. In PEnPF, parameters and states are updated simultaneously by EnKF
518 and PF, and the better estimations are shown as the posterior distributions. If at each time step,
519 EnKF performs better than PF, evolution characteristics of parameters and states would be
520 identical to those generated by EnKF. The results in [Figure 5](#) suggest that, parameter and state
521 estimations from EnKF are chosen as the corresponding posteriors in the data assimilation
522 experiment through PEnPF.

523

524 Moreover, to further explore the reliability of the four data assimilation approaches, five
525 sample size scenarios (i.e. {20, 50, 100, 200, 500}) are tested. For each scenario, the
526 synthetic experiment is performed for 30 replicates to identify the robustness of the proposed
527 approaches. The performances of EnKF, PF, CEnPF and PEnPF are evaluated through two
528 deterministic indices (i.e. RMSE and NSE) and two probabilistic indices (i.e. CRPS and

529 Resolution). [Figure 6](#) compares the performance of EnKF, PF, CEnPF and PEnPF through a
530 boxplot. The results indicate that all four methods will perform better with an increase in
531 sample size, and the sample size influence PF more significantly than the other three data
532 assimilation approaches. In detail, the PEnPF produce best deterministic predictions with
533 lowest values for NSE and RMSE, followed by EnKF, CEnPF and PF. The performance of
534 CEnPF is not as well as EnKF in this synthetic experiment. However, it performs better than
535 PF. Especially when the same size is larger than 50, CEnPF would generate more reliable
536 predictions than PF. For probabilistic predictions, the PEnPF would lead to lowest values for
537 CRPS, indicating closest distance between the predictive and observed cumulative
538 distribution functions (CDFs). Moreover, similar with deterministic predictions, the proposed
539 CEnPF does not perform as well as EnKF in this synthetic experiment, but it provide more
540 accurate predictions than PF, especially when the sample size is larger than 50.

541

542 **4. Real Case Study**

543 **4.1. Site Description**

544 Two real watersheds will be used test the applicability of the proposed data assimilation
545 schemes, as presented in [Figure 7](#). The first catchment is the Huanjiang river, located in the
546 northern part of Jing river basin with a drainage area of 4,640 km². This catchment has two
547 main tributaries, which converge together at Hongde (107.19 E, 36.76 N). In general, the Jing
548 river basin is characterized by a semi-arid and sub-humid continental monsoon climate,
549 resulting in significant temporal-spatial variations in precipitation. From the northern to
550 southern part, the corresponding annual precipitation ranges from 240 to 710 mm, with
551 approximately 50~60% precipitation occurring in the Summer and Fall seasons. In particular,
552 the Huanjiang in this case is located in the northern part of the Jing River watershed, and the

553 annual precipitation there fluctuates from 240 to 350 mm with mean annual precipitation of
554 approximate 309 mm. For Huanjiang river, the daily precipitation data from Ganjipan,
555 Fanxue, Shancheng, Wuqi, Gengwan, Honglaochi, Siheyuan and Hongde are employed to
556 generate areal precipitation over the entire sub-catchment. The potential precipitation values
557 were obtained through the Penman–Monteith equation, based on meteorological
558 measurements from national meteorological stations (i.e. Changwu, Xifengzhen, Guyuan,
559 Huanxian, Tongchuan) in the Jing river basin. [Tables 2 and 3](#) provide the location information
560 for the rain gauge stations and the national meteorological stations.

561

562 The second case is the Xiangxi river basin, located in the Three Gorges Reservoir area, China.
563 The Xiangxi river is located between 30.96 ~ 31.67 °N and 110.47 ~ 111.13°E in the Hubei
564 part of the China Three Gorges Reservoir (TGR) region, with a draining area of
565 approximately 3,200 km². The Xiangxi river originates in the Shennongjia Nature Reserve
566 with a main stream length of 94 km and a catchment area of 3,099 km² and is one of the main
567 tributaries of the Yangtze river ([Han et al., 2014](#); [Yang and Yang, 2014](#); [Miao et al., 2014](#)).
568 The watershed experiences a northern subtropical climate. The annual precipitation is about
569 1,100 mm and ranges from 670 to 1,700 mm with considerable spatial and temporal
570 variability ([Xu et al., 2010](#); [Zhang et al., 2014](#)). The main rainfall season is from May
571 through September, with a flooding season from July to August. The annual average
572 temperature in this region is 15.6 °C and ranges from 12 °C to 20 °C. For this case,
573 meteorological and streamflow data at Xingshan (31°13'N, 110°45'E) station will be used.

574

575 -----

576 Place Figure 7 here and Tables 2 and 3 here

577 -----

578 **4.2. Results Analysis for Huanjiang river**

579 In hydrologic sequential data assimilation, two issues are generally predefined before
580 implementation of the sequential data assimilation. The first one is how many ensembles or
581 particles are going to use to represent the distributional information in parameters, state
582 variables and predictions. The other one is that how to account for uncertainty existing in
583 forcing data, model prediction, and streamflow measurements. In the real case study, the
584 sample size is set to be 200 for all the four data assimilation schemes based on the results of
585 the synthetic experiment. Moreover, random perturbations are added to model inputs, outputs,
586 and parameters to reflect their inherent uncertainties. In this study, the precipitation is
587 assumed to follow a lognormal distribution with the proportional factors being 20% of the
588 true, while the potential evapotranspiration, streamflow measurements, and model prediction
589 are normally distributed with the standard errors being 20% of the true values.

590

591 Figure 8 shows the comparison between ensemble predictions of the four data assimilation
592 methods and observations. Figure 8(a) indicates the comparison between the mean
593 predictions and predictive intervals from EnKF and model and observations. The result
594 shows that the predictive intervals from EnKF can generally bracket the observations during
595 the low flow period, while underestimations occur during the high flow period. Similar
596 characteristics can be found for both PF. However, as shown in Figure (8b), PF provide better

597 predictions than EnKF. Especially for the high flow periods, the predictive intervals from PF
598 can catch the peak flow better than those from EnKF. In comparison with EnKF and PF, the
599 proposed CEnPF can generate more reliable predictions. As shown in Figure (8c), the
600 predictive intervals from CEnPF can generally bracket the observations while the ensemble
601 means can well track the fluctuation of real discharges for both low and high flow periods.
602 For the PEnPF, it seems to perform slightly worse than CEnPF. In particular, the PEnPF
603 would generate worse (i.e. underestimation) predictions than PF during the high flow periods.
604 However, the PF would produce overestimations in a quite long period after the highest peak
605 flow while PEnPF can provide accurate predictions in this period. In this case, the predictions
606 from CEnPF lead to a NSE value of 0.911, a RMSE value of 5.897, a CRPS value of 2.209
607 and a Resolution value of 41.685. The four indices (i.e. NSE, RMSE, CRPS and Resolution)
608 correspond to the predictions of PEnPF are 0.861, 7.372 , 1.675 and 15.058, respectively. The
609 four indices for the predictions of EnKF are 0.767, 9.540, 2.234, and 21.697, and those
610 indices for PF predictions are 0.776, 9.354, 4.026, and 38.596. Consequently, the CEnPF
611 leads to best deterministic predictions while the PEnPF generates best probabilistic
612 predictions

613

614 -----

615 Place Figure 8 here

616 -----

617

618 To further demonstrate the applicability of the proposed data assimilation methods, four
619 sample scenarios (i.e. {50, 100, 200, 500}) are further tested for this real case with 10
620 replicates conducted for each sample scenario. [Figure 9](#) compares the performance of EnKF,
621 PF, CEnPF and PEnPF through a boxplot. It shows that as the increase of sample size, the
622 proposed CEnPF, PEnPF as well as traditional EnKF would generate reliable predictions with
623 the four evaluation indices varied within limited intervals. In comparison, the PF can also
624 generate unsatisfactory results even the sample size of 500. Tables 4 to 7 provide the mean,
625 minimum and maximum values for NSE, RMSE, CRPS and Resolution for the 10 replicates
626 by different data assimilation schemes under different sample size scenarios. The results
627 indicate that the proposed CEnPF can generally provide best results for deterministic
628 predictions with lowest NSE and RMSE values. For instance, the CEnPF can lead to a mean
629 NSE value of 0.78 under a sample size of 100, which is higher than the other three
630 approaches (i.e. the mean NSE values would be 0.72, 0.69 and 0.65 for PEnPF, EnKF and
631 PF). In comparison, the PEnPF would produce better probabilistic predictions than CEnPF,
632 EnKF and PF, which generally has lowest CRPS and Resolution values, as presented in
633 Tables 6 and 7. In general, even though the prediction from CEnPF has large degree of
634 uncertainty (i.e. large Resolution values), the proposed CEnPF and PEnPF can provide better
635 results for both deterministic and probabilistic forecasts for the Huanjiang river basin

636

637 -----

638 Place Figure 9, Tables 4 to 7 here

639 -----

640

641 **4.3. Results Analysis for Xiangxi river**

642

643 The developed data assimilation approaches are further applied for hydrological data
644 assimilation in Xiangxi river, which is an main tributary of Yangtze river in Hubei Province.
645 The Xiangxi river basin experiences a northern subtropical climate with higher temperature
646 and precipitation than the Huanjiang river basin which has a semi-arid climate. To clearly
647 account uncertainties in meteorological data and streamflow measurements in Xiangxi river,
648 the proportional factor is set to be 30% of the true measurements. In current case, the sample
649 size is 500.

650

651 Figure 10 shows the performance of the developed CEnPF and PEnPF as well as traditional
652 EnKF and PF approaches for hydrological data assimilation in Xiangxi river. As presented in
653 Figure (10a), the EnKF approach provide accurate deterministic and probabilistic predictions
654 during the low flow periods, but these predictions cannot well track observations during high
655 flow periods and show underestimated results in these periods. Compared with EnKF, the PF
656 approach seems to provide better predictions, as shown in Figure (10b). Especially in high
657 flow periods, PF performs better than EnKF, but it still provides underestimations in these
658 time steps. In comparison, the developed CEnPF and PEnPF are able to generate reliable
659 results for both deterministic predictions and the associated predictive intervals. As shown in
660 Figures (10c) and (10), the predictive intervals of CEnPF and PEnPF can bracket the real
661 observations at most time periods for this case. Meanwhile, the corresponding deterministic

662 predictions (i.e. predictive means) can trace the variation in streamflow in both high and low
663 flow periods.

664 -----

665 Place Figure 10

666 -----

667

668 Table 8 shows the performance of the four approaches for hydrological data assimilation in

669 Xiangxi river basin under different sample size scenarios. The results shows that for

670 deterministic predictions, the proposed CEnPF and PEnPF approach performs better than

671 EnKF in all selected sample scenarios, and these two methods provide better deterministic

672 predictions than PF in three of the four sample scenarios. However, in terms of the

673 probabilistic forecasts, the performances of the fours approaches show different features.

674 EnKF seems to lead to lowest CRPS values for all sample scenarios. However, at least one

675 proposed approach (i.e. CEnPF or PEnPF) can provide better probabilistic predictions than

676 PF for all selected sample scenarios.

677 -----

678 Place Tables 8 here

679 -----

680

681 **5. Discussion**

682 In this study, both CEnPF and PEnPF integrate traditional PF and EnKF into combined

683 framework. This means that the computational demand would increase for CEnPF and

684 PEnPF since they have additional procedures. Figure 11 presents the computation demand for
685 EnKF, PF, CEnPF and PEnPF. The results show that, among these four approaches, PF
686 requires least computational time, and both CEnPF and PEnPF require more computational
687 time than EnKF and PF since they have more steps. However, the computational time for the
688 two developed methods is manageable. In detail, the PEnPF needs more computational
689 requirement than the other three approaches. For instance, the computational time for PEnPF
690 would be about 590 seconds when the sample size is 500, while the time for EnKF, PF and
691 PEnPF would be 347, 102 and 443 seconds, respectively. This is because that, in spite of
692 update procedures of EnKF and PF, the PEnPF needs two additional steps for putting the
693 updated parameters from EnKF and PF into the original hydrological model to evaluate the
694 mismatch between the resulting outputs and the real observations at each time step. This
695 suggests that for some large hydrological models requiring much computation time, the
696 PEnPF may need much more time than EnKF, PF and PEnPF since the hydrological model
697 would be run for $3*ns$ (ns is the sample size) times at each time while the other three
698 approaches only need to run the hydrological model ns times.

699 -----

700 Place Figure 11 here

701 -----

702 **6. Conclusions**

703 This study proposed two integrated data assimilation schemes, i.e. the coupled EnKF and PF
704 (CEnPF) and the parallelized EnKF and PF (PEnPF) approaches through the integration of
705 the capabilities of EnKF and PF. The CEnPF sequentially adopts EnKF and PF to update

706 model parameters and states, in which EnKF is first applied to correct model states and
707 parameters, and PF is then employed to eliminate insignificant particles. In comparison, the
708 PEnPF approach simultaneously updates model states and parameters in parallel through
709 EnKF and PF, and chooses the better estimates as the posterior distributions. The proposed
710 CEnPF and PEnPF approaches were applied for hydrologic data assimilation in two
711 real-world cases to demonstrate their applicability in quantifying uncertainty in hydrologic
712 prediction

713

714 A synthetic application firstly illustrated procedures of the proposed CEnPF and PEnPF
715 approaches and compared them with traditional PF and EnKF methods. Five sample size
716 scenarios were tested to evaluate the performance of the proposed methods. The results
717 suggested that PEnPF performed best for both probabilistic and deterministic predictions,
718 while CEnPF could provide better predictions than PF. The improvement of the proposed
719 CEnPF and PEnPF upon EnKF and PF was further illustrated by two real-world catchments
720 with different climate conditions. The results for the Huanjiang river, located in the northern
721 part of Jing river, demonstrated that PEnPF would produce better probabilistic predictions
722 than CEnPF, EnKF and PF, which generally has lowest CRPS and Resolution and the CEnPF
723 could provide better results in deterministic predictions but lead to large uncertainty in its
724 ensemble outputs. For the Xiangxi river located in the Yangtze river basin, the results
725 indicated that the proposed approach improved EnKF and PF in terms of deterministic
726 predictions. For all selected sample size scenarios, at least one method could give better
727 probabilistic predictions than PF.

728

729 The ensemble Kalman filter (EnKF) and particle filter (PF) methods have been extensively
730 applied for hydrologic data assimilation. However, both of them have their inherent
731 disadvantages which restrict their application for many cases. In this study, two integrated
732 sequential data assimilation approaches are proposed by integrating the capabilities of EnKF
733 and PF into a general framework. The case studies for synthetic experiment and two
734 real-world hydrologic data assimilation problems demonstrate the significant potential of the
735 proposed CEnPF and PEnPF approaches. Moreover, the computational time for CEnPF and
736 PEnPF is manageable when compared with EnKF and PF. However, the PEnPF may require
737 much more computational time for large-scale or time-consuming hydrological models than
738 EnKF, PF and CEnPF.

739

740

741 **Acknowledgement**

742 This work was jointly funded by the Natural Science Foundation of China (51520105013),
743 the National Key Research and Development Plan (2016YFC0502800), and the Natural
744 Sciences and Engineering Research Council of Canada.

745 **References**

- 746 Ajami N.K., Duan Q.Y., Sorooshian S., (2007). An integrated hydrologic Bayesian multimodel
747 combination framework: Confronting input, parameter, and model structural uncertainty in hydrologic
748 prediction. *Water Resources Research*, 43, W01403.
- 749 Brown, J. D. (2010). Prospects for the open treatment of uncertainty in environmental research, *Prog. Phys.*
750 *Geog.*, 34, 75–100, doi:10.1177/0309133309357000.
- 751 Chen, H., Yang, D., Hong, Y., Gourley, J.J., Zhang, Y., 2013. Hydrological data assimilation with the
752 Ensemble Square-Root-Filter: use of streamflow observations to update model states for real-time flash
753 flood forecasting. *Advance in Water Resources* 59, 209-220.
- 754 Chen Y., Oliver D.S., (2013). Levenberg–Marquardt forms of the iterative ensemble smoother for efficient
755 history matching and uncertainty quantification. *Computational Geosciences* 17(4), 689-703.
- 756 Clark, M. P., Rupp, D. E., Woods, R. A., Zheng, X., Ibbitt, R. P., Slater, A. G., Schmidt, J., and Uddstrom,
757 M. J., (2008). Hydrological data assimilation with the ensemble Kalman filter: Use of streamflow
758 observations to update states in a distributed hydrological model, *Adv. Water Resour.*, 31, 1309–1324,
759 doi:10.1016/j.advwatres.2008.06.005.
- 760 DeChant C.M., Moradkhani H., (2012). Examining the effectiveness and robustness of sequential data
761 assimilation methods for quantification of uncertainty in hydrologic forecasting. *Water Resources*
762 *Research*, 48, W04518, doi:10.1029/2011WR011011
- 763 DeChant C.M., and H. Moradkhani (2014), Toward a Reliable Prediction of Seasonal Forecast Uncertainty:
764 Addressing Model and Initial Condition Uncertainty with Ensemble Data Assimilation and Sequential
765 Bayesian Combination, *Journal of Hydrology* , 519, 2967-2977, doi: 10.1016/j.jhydrol.2014.05.045.
- 766 De Lannoy, G. J. M., Reichle, R. H., Houser, P. R., Pauwels, V. R., and Verhoest, N. E., (2007). Correcting
767 for forecast bias in soil moisture assimilation with the ensemble Kalman filter, *Water Resour. Res.*, 43,
768 W09410, doi:10.1029/2006WR005449.
- 769 Doucet, A., N. De Freitas, and N. Gordon (2001), *Sequential Monte Carlo Methods in Practice*, vol. 1,
770 Springer, N. Y.
- 771 Evensen, G. (1994). Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte
772 Carlo methods to forecast error statistics. *Journal of Geophysical Research: Oceans*, 99(C5),
773 10143-10162.

774 Evensen, G. (2003), The Ensemble Kalman Filter: theoretical formulation and practical implementation,
775 Ocean Dynamics 53, 343–367.

776 Fan Y.R., Huang G.H., Baetz B.W., Li Y.P., Huang K., Li Z., Chen X., Xiong L.H., (2016). Parameter
777 uncertainty and temporal dynamics of sensitivity for hydrologic models: A hybrid sequential data
778 assimilation and probabilistic collocation method. *Environmental Modelling & Software* 86, 30-49

779 Fan Y.R., Huang G.H., Huang K., Baetz B.W., (2015a). Planning Water Resources Allocation under
780 Multiple Uncertainties through A Generalized Fuzzy Two-Stage Stochastic Programming Method.
781 IEEE Transactions on Fuzzy Systems, 23(5), 1488-1504.

782 Fan Y.R., Huang W.W., Li Y.P., Huang G.H., Huang K., Li Y.P., (2015b). A coupled ensemble filtering and
783 probabilistic collocation approach for uncertainty quantification of hydrological models. *Journal of*
784 *Hydrology*, 530, 255-272.

785 Fan Y.R., Huang W.W., Huang G.H., Huang K., Zhou X., (2015c). A PCM-based stochastic hydrological
786 model for uncertainty quantification in watershed systems. *Stochastic Environmental Research and*
787 *Risk Assessment*, 29, 915-927

788 Frei, M., and H. R. Kunsch (2013), Bridging the ensemble Kalman and particle filters, *Biometrika*, 100(4),
789 781–800

790 Fearnhead, P. and Clifford, P.: On-line inference for hidden Markov models via particle filters, *J. R. Stat.*
791 *Soc. B Met.*, 65, 887–899, 2003.

792 Gordon, N.J., Salmond, D.J., Smith, A.F.M., 1993. Novel approach to nonlinear/nonGaussian Bayesian
793 state estimation. *IEEE Proceedings F: Radar Signal Process.* 140 (2), 107e113.

794 Gu Y., Oliver D.S., (2007). An iterative ensemble Kalman filter for multiphase fluid flow data assimilation.
795 *SPE Journal*, 12(4), 438–446

796 Han X., Li X., (2008). An evaluation of the nonlinear/non-Gaussian filters for the sequential data
797 assimilation. *Remote Sensing of Environment*, 112, 1434-1449.

798 Kavetski, D., Kuczera, G., and Franks, S. W. (2006a). Bayesian analysis of input uncertainty in
799 hydrological modeling: 1. Theory, *Water Resour. Res.*, 42, W03407, doi:10.1029/2005WR00436.

800 Kavetski, D., Kuczera, G., and Franks, S. W., (2006b) Bayesian analysis of input uncertainty in
801 hydrological modeling: 2. Application, *Water Resour. Res.*, 42, W03408, doi:10.1029/2005WR004376.

802 Kong X.M., Huang G.H., Fan Y.R., Li Y.P., (2015). Maximum entropy-Gumbel-Hougaard copula method

803 for simulation of monthly streamflow in Xiangxi river, China. *Stochastic Environmental Research and*
804 *Risk Assessment* 29, 833-846.

805 Leisenring, M., & Moradkhani, H. (2012). Analyzing the uncertainty of suspended sediment load
806 prediction using sequential data assimilation. *Journal of Hydrology*, 468, 268-282

807 Li T., Bolic M., Djuric P.M., (2015) Resampling Methods for Particle Filtering: Classification,
808 implementation, and strategies. *IEEE Signal Processing Magazine*, 32(3), 70-86

809 Li Z., Huang G.H., Fan Y.R., Xu J.L., Hydrologic Risk Analysis for Nonstationary Streamflow Records
810 under Uncertainty. *Journal of Environmental Informatics* 26 (1), 41-51.

811 Liu, J. S. and Chen, R.: Sequential Monte Carlo methods for dynamic systems, *J. Am. Stat. Assoc.*, 93,
812 1032–1044, 1998.

813 Liu Y., Weerts A.H., Clark M., Hendricks Franssen H.-J., Kumar S., Moradkhani H., Seo D.-J.,
814 Schwanenberg D., Smith P., van Dijk A.I.J.M., van Velzen N., He M., Lee H., Noh S.J., Rakovec O.,
815 Restrepo P., (2012). Advancing data assimilation in operational hydrologic forecasting: progresses,
816 challenges, and emerging opportunities. *Hydrology and Earth System Sciences*, 16, 3863-3887.

817 Madadgar, S. and H. Moradkhani (2014). Improved Bayesian Multi-modeling: Integration of Copulas and
818 Bayesian Model Averaging. *Water Resources Research*, 50, 9586-9603, doi: 10.1002/2014WR015965.

819 Montanari A., Brath A., (2004), A stochastic approach for assessing the uncertainty of rainfall-runoff
820 simulations. *Water Resources Research*, 40, W01106

821 Moor, R.J., 1985. The probability-distributed principle and runoff production at point and basin scales.
822 *Hydrological Science Journal* 30, 273-297.

823 Moor, R.J., 2007. The PDM rainfall-runoff model. *Hydrology and Earth Systems Science* 11 (1), 483-499.

824 Moradkhani, H. (2008). Hydrologic remote sensing and land surface data assimilation. *Sensors*, 8(5),
825 2986-3004.

826 Moradkhani, H., S. Sorooshian, H. V. Gupta, and P. Houser (2005a), Dual state – parameter estimation of
827 hydrologic models using ensemble Kalman filter. *Advances in Water Resources*, 28, 135 – 147.

828 Moradkhani H., Dechant C.M., Sorooshian S., (2012). Evolution of ensemble data assimilation for
829 uncertainty quantification using the particle filter-Markov chain Monte Carlo method, *Water Resources*
830 *Research*, 48, W12520, doi:10.1029/2012WR012144.

831 Moradkhani, H., Hsu, K. L., Gupta, H., & Sorooshian, S. (2005b). Uncertainty assessment of hydrologic

832 model states and parameters: Sequential data assimilation using the particle filter. *Water Resources*
833 *Research*, 41(5)

834 Pappenberger, F. and Beven, K. J., (2006). Ignorance is bliss: Or seven reasons not to use uncertainty
835 analysis, *Water Resour. Res.*, 42, W05302, doi:10.1029/2005WR004820, 2006.

836 Parrish, M., H. Moradkhani, and C.M. DeChant (2012). Towards Reduction of Model Uncertainty:
837 Integration of Bayesian Model Averaging and Data Assimilation, *Water Resources Research*, 48,
838 W03519, doi:10.1029/2011WR011116.

839 Pathiraja, S., L. Marshall, A. Sharma, and H. Moradkhani (2016), Detecting non-stationary hydrologic
840 model parameters in a paired catchment system using Data Assimilation, *Advances in Water Resources*,
841 94, 103-119, doi:10.1016/j.advwatres.2016.04.021.

842 Pathiraja, S., L. Marshall, A. Sharma, and H. Moradkhani (2016), Hydrologic Modeling in Dynamic
843 catchments: A Data Assimilation Approach, *Water Resources Research*, doi: 10.1002/2015WR017192

844 Plaza-Guingla D. A., De Keyser R., De Lannoy G. J. M., Giustarini L., Matgen P., and Pauwels V. R. N.,
845 (2013). Improving particle filters in rainfall-runoff models: Application of the resample-move step
846 and the ensemble Gaussian particle filter, *Water Resource Research*, 49, doi:10.1002/wrcr.20291.

847 Rasmussen J., Madsen H., Jensen K.H., Refsgaard J.C., (2015). Data assimilation in integrated
848 hydrological modeling using ensemble Kalman filtering: evaluating the effect of ensemble size and
849 localization on filter performance. *Hydrology and Earth System Sciences*, 19, 2999-3013.

850 Reichle R., McLaughlin D., Entekhabi D., (2002). Hydrologic data assimilation with the ensemble Kalman
851 filter. *Monthly Weather Review*, 130(1), 103-114.

852 Rezaie, J. and Eidsvik, J. (2012). Shrunked $(1 - \alpha)$ ensemble Kalman filter and α Gaussian mixture filter.
853 *Computational Geosciences*, 16:837–852.

854 Schaake, J., Franz, K., Bradley, A., and Buizza, R., (2006). The Hydrologic Ensemble Prediction
855 EXperiment (HEPEX), *Hydrol. Earth Syst. Sci. Discuss.*, 3, 3321–3332,
856 doi:10.5194/hessd-3-3321-2006.

857 Salamon, P. and Feyen, L. (2010). Disentangling uncertainties in distributed hydrological modeling using
858 multiplicative error models and sequential data assimilation, *Water Resour. Res.*, 46, W12501,
859 doi:10.1029/2009WR009022.

860 Shen Z., and Tang Y., (2015). A modified ensemble Kalman particle filter for non-Gaussian systems with

861 nonlinear measurement functions. *Journal of Advances in Modeling Earth Systems*, 07,
862 doi:10.1002/2014MS000373.

863 Shi, Y., K. J. Davis, F. Zhang, C. J. Duffy, and X. Yu (2014), Parameter estimation of a physically based
864 land surface hydrologic model using the ensemble Kalman filter: A synthetic experiment. *Water*
865 *Resources Research*, 50, 1-19, doi:10.1002/2013WR014070

866 Snyder, C., Bengtsson, T, Bickel, P., and Anderson, J.: Obstacles to high-dimensional particle filtering,
867 *Mon. Weather Rev.*, 136, 4629–4640, 2008.

868 Su, C.H., Ryu, D., Crow, W.T., Western, A.W., 2014. Beyond triple collocation: applications to soil
869 moisture monitoring. *J. Geophys. Res. Atmos.* 119, 6419–6439.
870 <http://dx.doi.org/10.1002/2013JD021043>.

871 Vrugt, Jasper A., ter Braak, Cajo J.F., Diks Cees G.H., Schoups, Gerrit, (2013). Hydrologic data
872 assimilation using particle Markov chain Monte Carlo simulation: Theory, concepts and applications.
873 *Advances in Water Resources*, 51, 457-478.

874 Vrugt J.A., Diks C.G.H., Gupta H.V., Bouten W., Verstraten J.M., (2005). Improved treatment of
875 uncertainty in hydrologic modelling: Combining the strengths of global optimization and data
876 assimilation. *Water Resources Research*, 41, W01017.

877 Weerts A.H., El Serafy G.Y.H., (2006). Particle filtering and ensemble Kalman filtering for state updating
878 with hydrological conceptual rainfall-runoff models. *Water Resources Research*, 42, W09403.
879 doi:10.1029/2005WR004093

880 Xie X., Zhang D., (2013). A partitioned update scheme for state-parameter estimation of distributed
881 hydrologic models based on the ensemble Kalman filter. *Water Resources Research*, 49, 7530-7365

882 Yan, H., DeChant, C.M., Moradkhani, H., 2015. Improving soil moisture profile prediction with the
883 particle filter-Markov chain Monte Carlo method. *IEEE Transactions on Geoscience and Remote*
884 *Sensing*. 53, 6134–6147. <http://dx.doi.org/10.1109/TGRS.2015.2432067>.

885 Yan, H., C.M., Moradkhani, H., 2016. Combined assimilation of streamflow and satellite soil moisture
886 with the particle filter and geostatistical modeling. *Advances in Water Resources*, 94, 364–378

887 Zhang Y., Oliver D.S., Chen Y., Skaug H.J., (2014). Data Assimilation by Use of the Iterative Ensemble
888 Smoother for 2D Facies Models, *SPE Journal*, 20(1), 169–185.

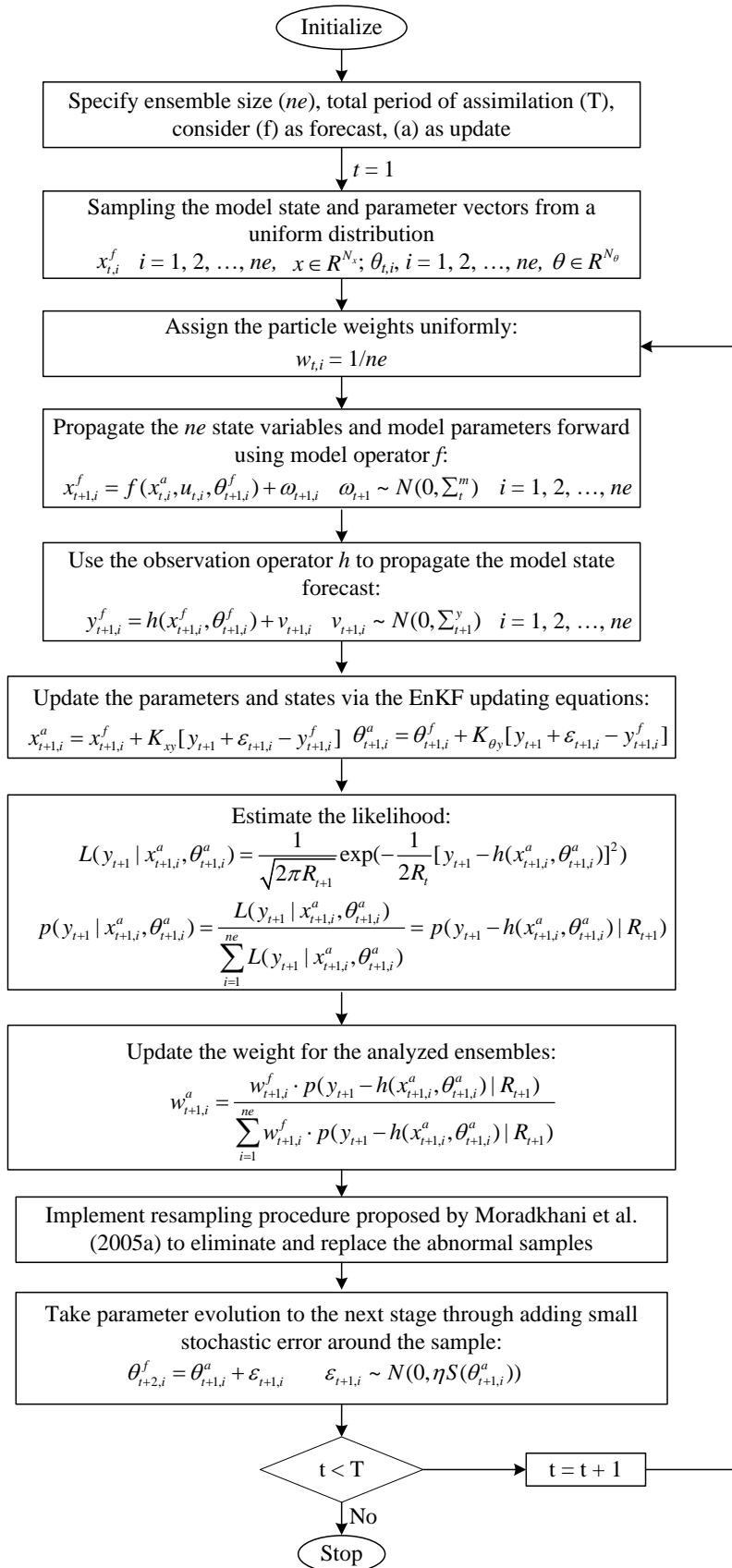


Figure 1. The flow chart of CEnPF

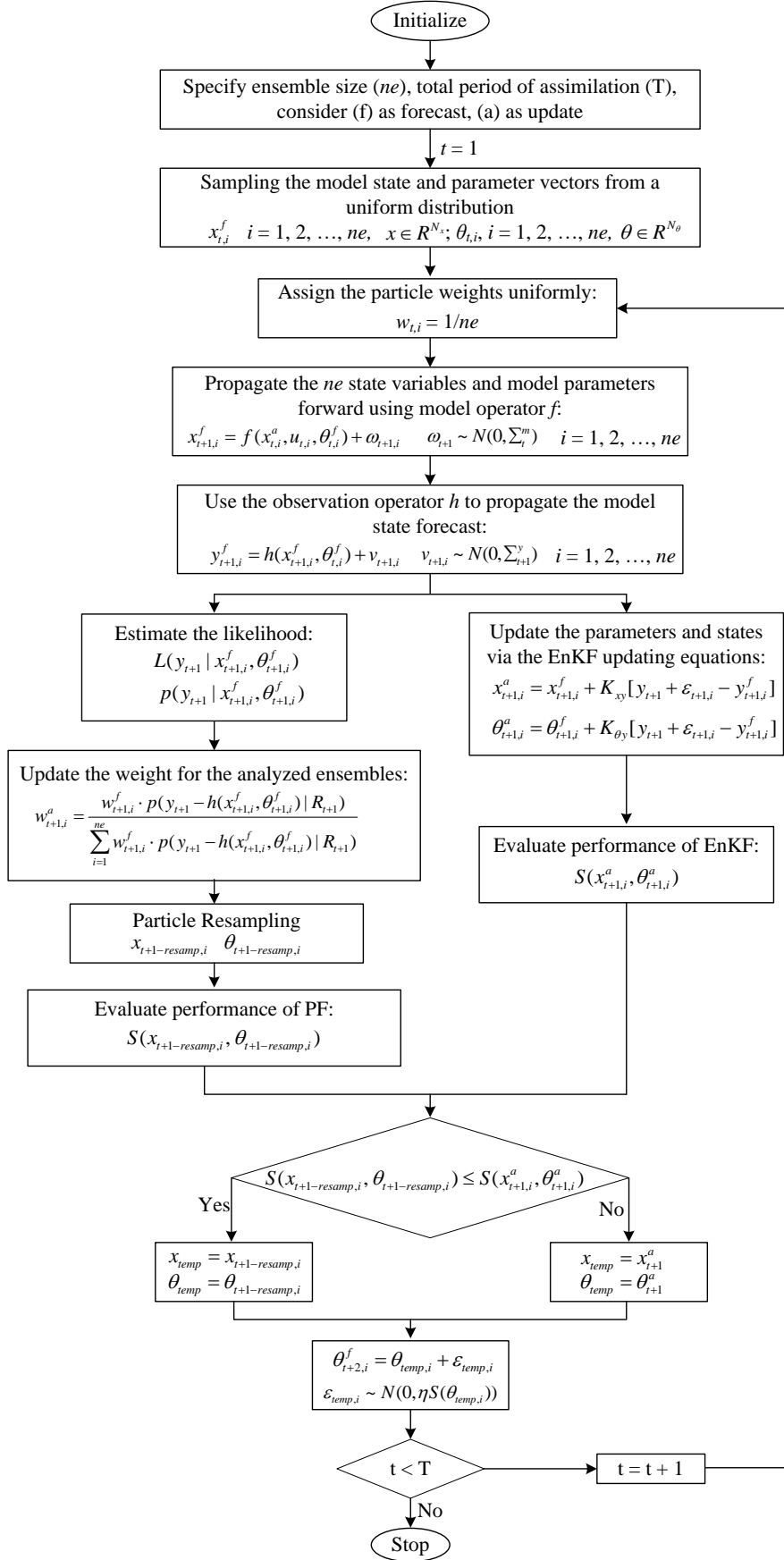
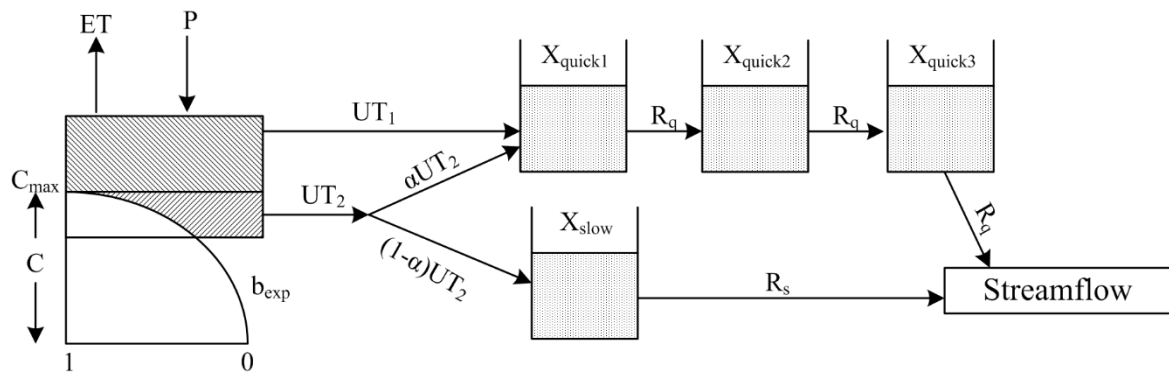


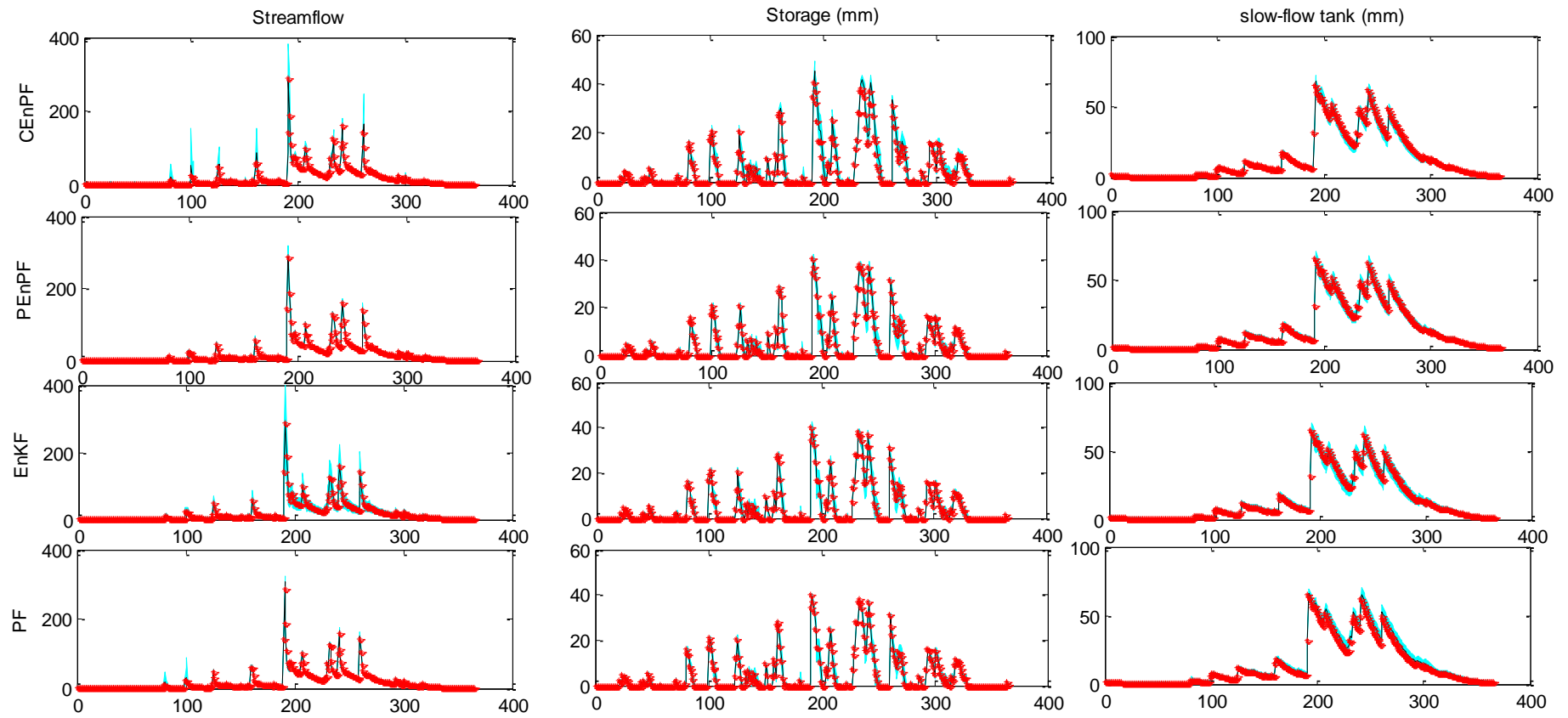
Figure 2. The flow chart of PENPF

894



895

896 Figure 3 Description of Hymod



897

898 Figure 4: Comparison between ensemble predictions and synthetically generated true discharge: Four methods are used including EnKF, PF, CEnPF and PEnPF. The cyan
 899 region indicates the 90% predictive intervals, the red stars denote the synthetic observations, and the black line indicates the predictive mean values.

900

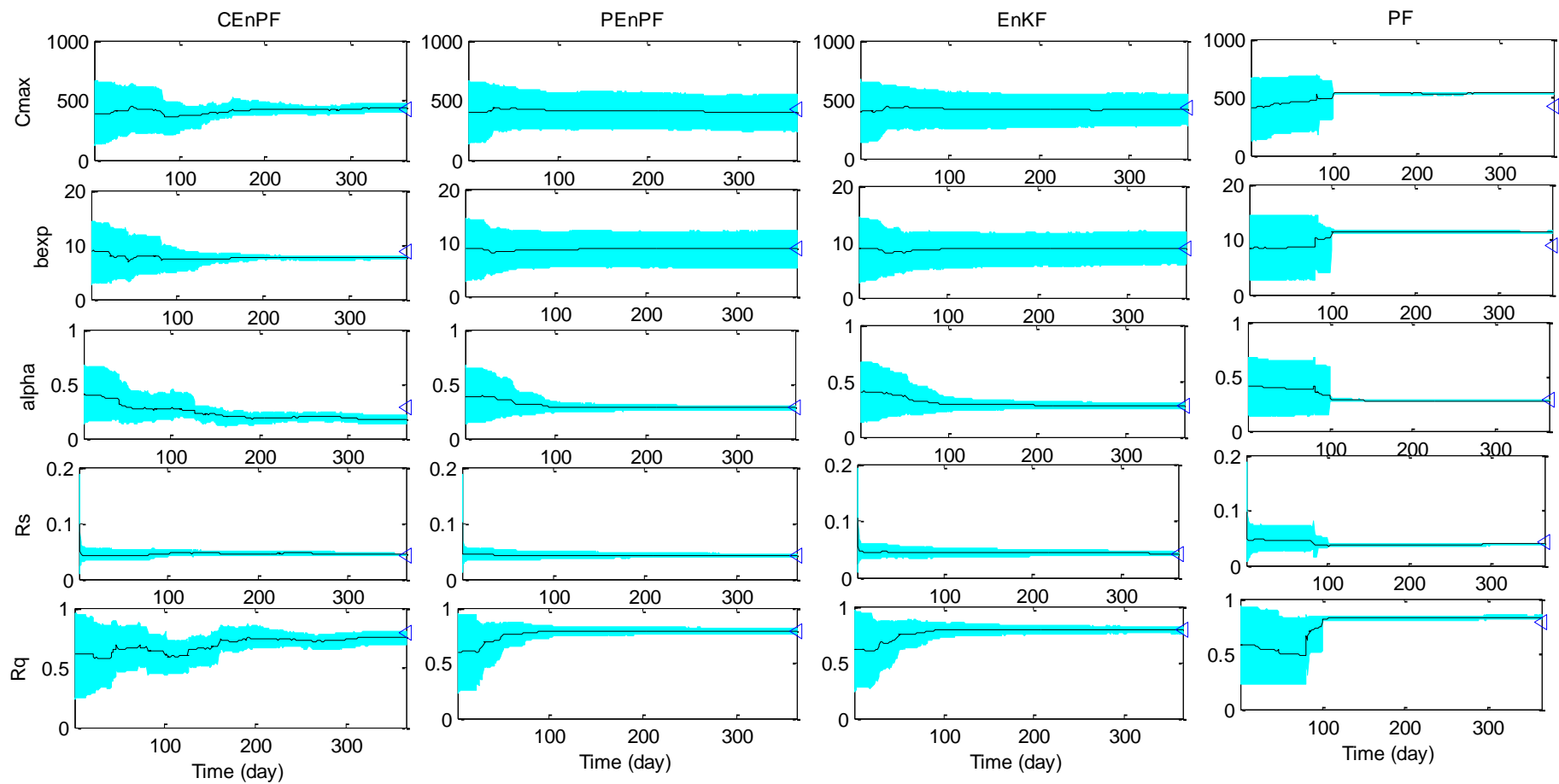


Figure 5: Convergence of the parameter distributions for the EnKF, PF, CEnPF and PEnPF for the synthetic experiments: The cyan region indicates the 90% intervals, the black line denotes the mean values, and the triangle is the predefined parameter value.

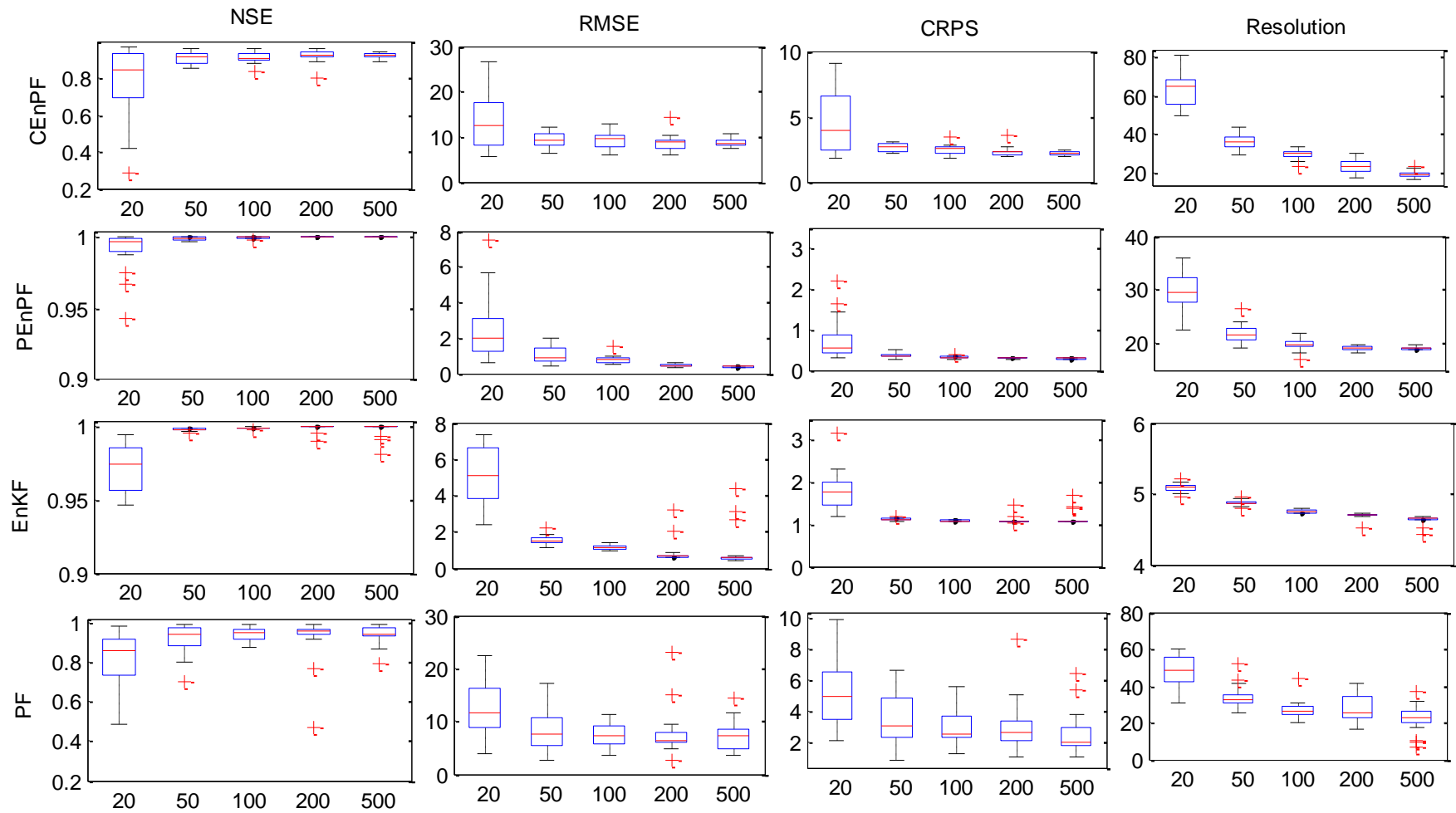


Figure 6. Performance comparison among EnKF, PF, CEnPF and PEnPF through a boxplot: The results show that all four methods will perform better with an increase in sample size. Generally, the PEnPF performs best than the other in both deterministic and probabilistic predictions, followed by EnKF, CEnPF and PF, if they are evaluated through NSE, RMSE and CRPS. However, the EnKF produces predictions with a lower resolution than PEnPF.

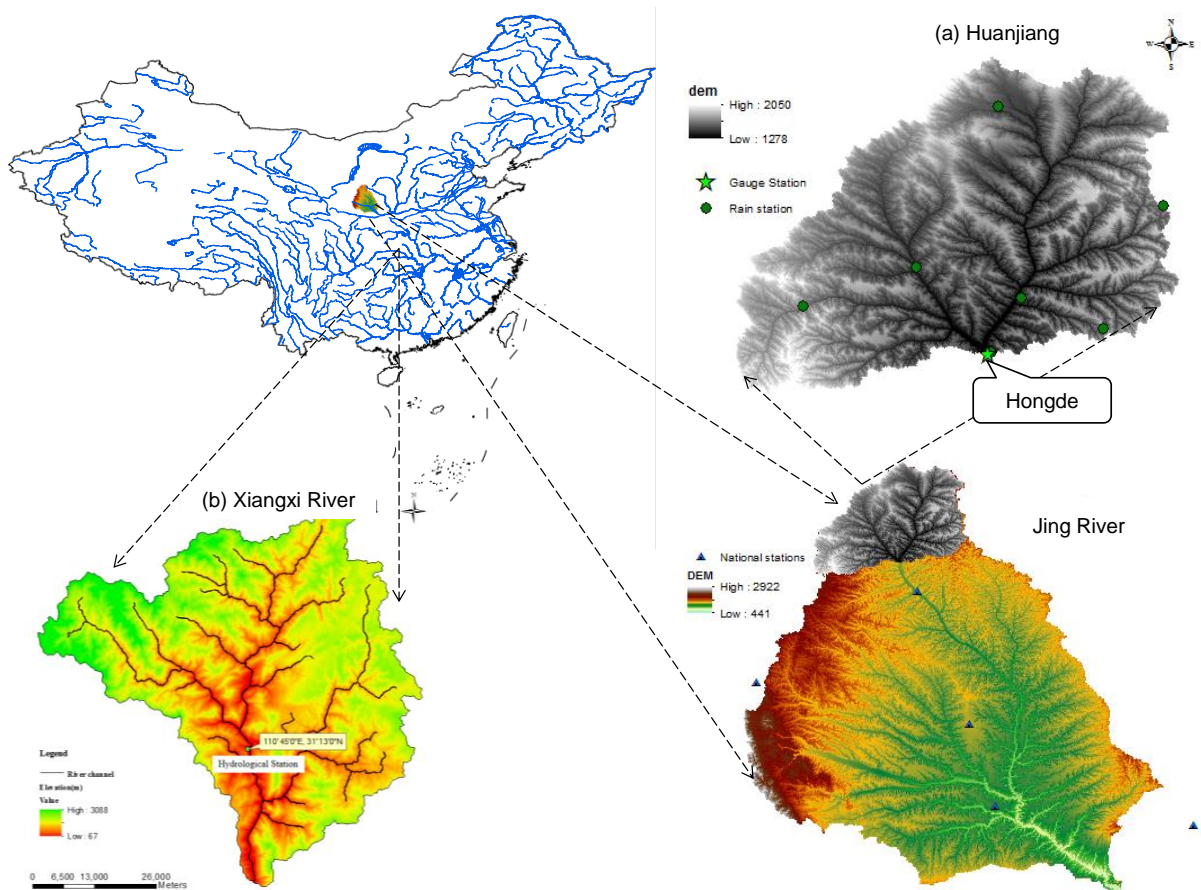


Figure 7. The location of the studied watersheds. Two watersheds are used to demonstrate the applicability of the proposed data assimilation schemes. One watershed named Huanjiang, located in the the north part of Jing River. Precipitation data from the seven rain stations in this catchment are used to generate the areal precipitation in the studied sub-catchment. The potential evapotranspiration (PE) are interpolated based on the PE results at the five national meteorological stations. The streamflow observations at Hongde station are used to evaluate the performance of the proposed methods. For the Xiangxi river watershed, meteorological and streamflow observations at Xingshan ($31^{\circ}13'N$, $110^{\circ}45'E$) station will be used.

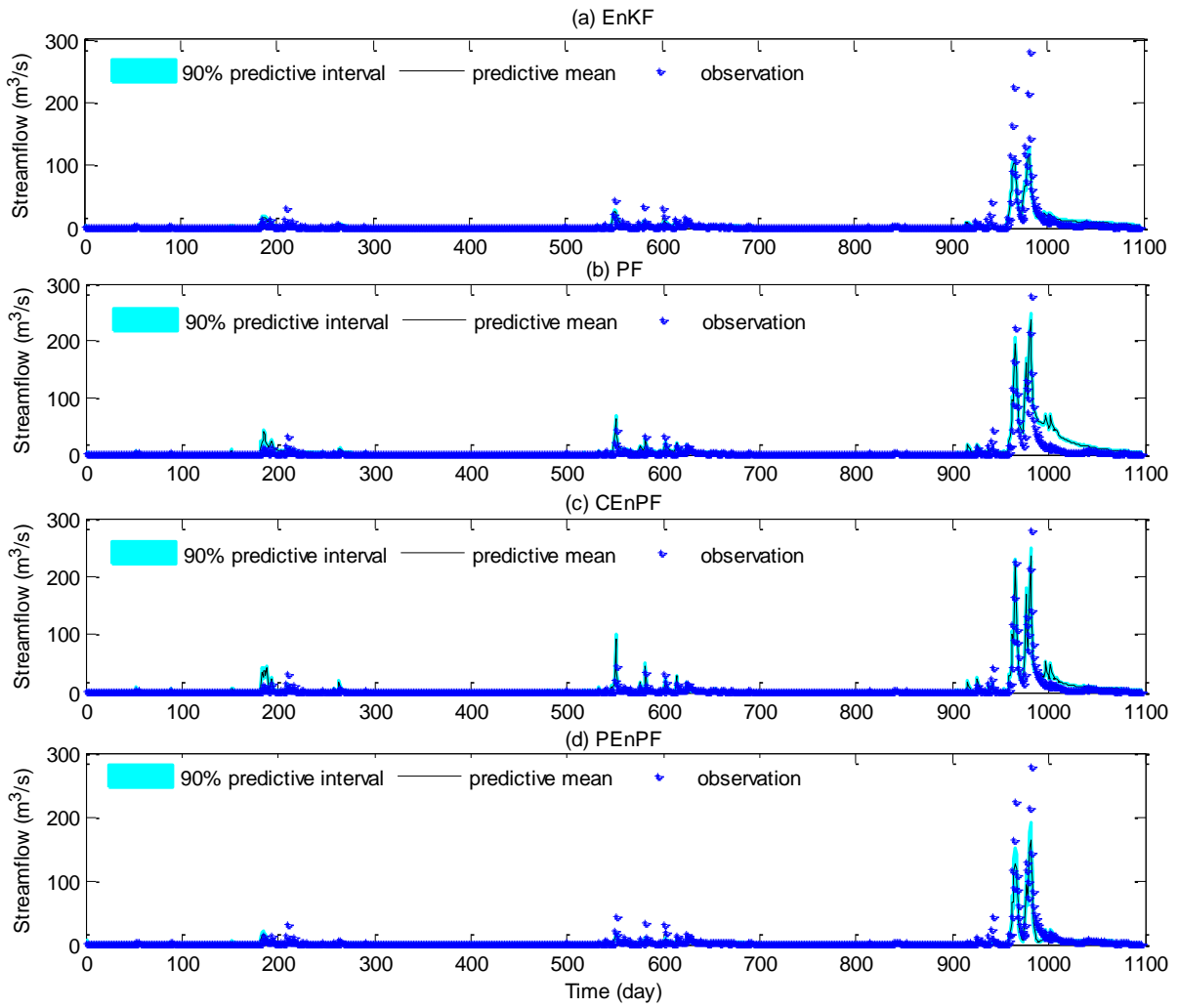


Figure 8. Comparison between the prediction intervals and observations for Huanjiang river through different data assimilation schemes: (a) EnKF, (b) PF, (c) CEnPF, (d) PEnPF.

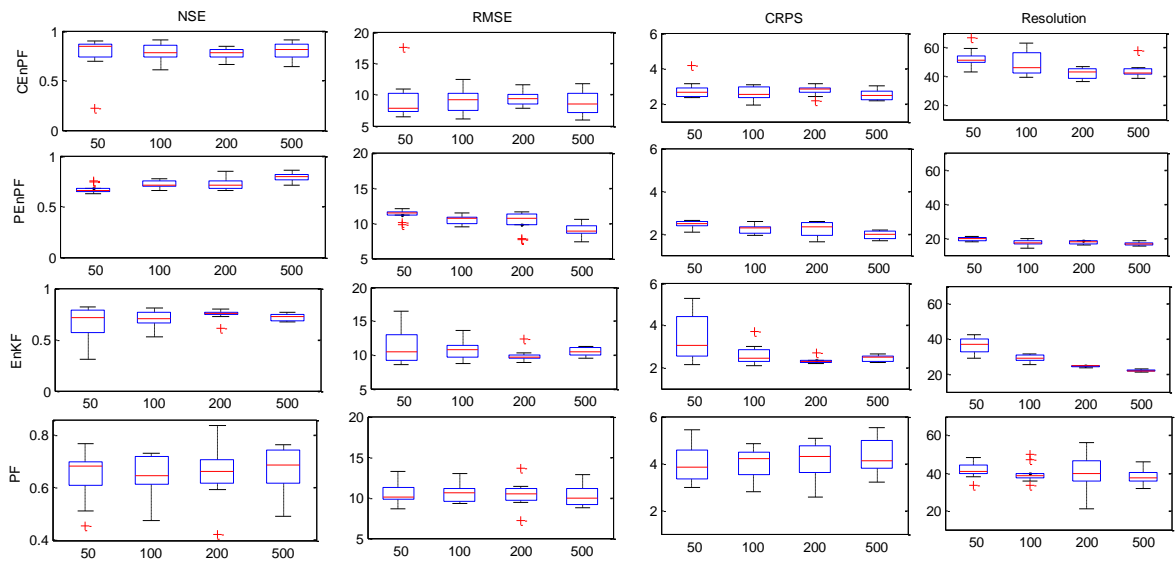


Figure 9. Performance comparison among different data assimilation schemes by using NSE, RMSE, CRPS and Resolution

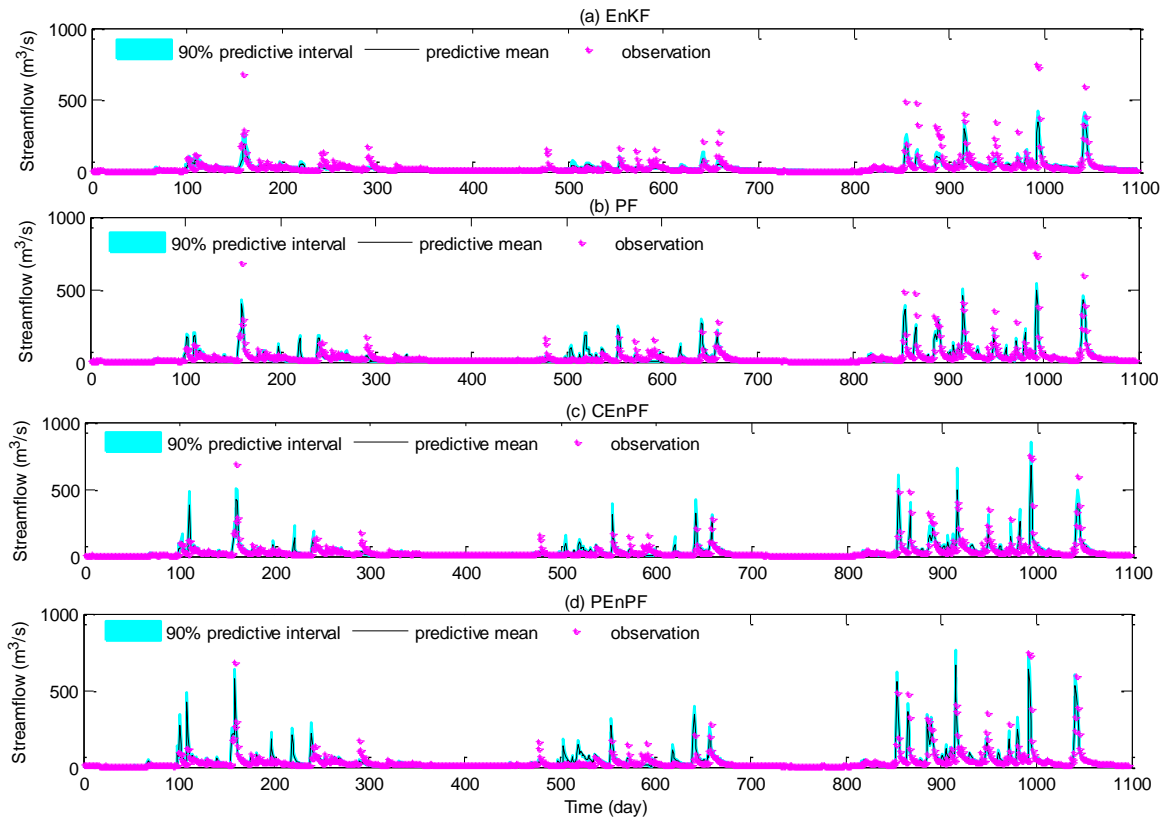


Figure 10. Comparison between the predication intervals and observations for Xiangxi river through different data assimilation schemes: (a) EnKF, (b) PF, (c) CEnPF, (d) PEnPF.

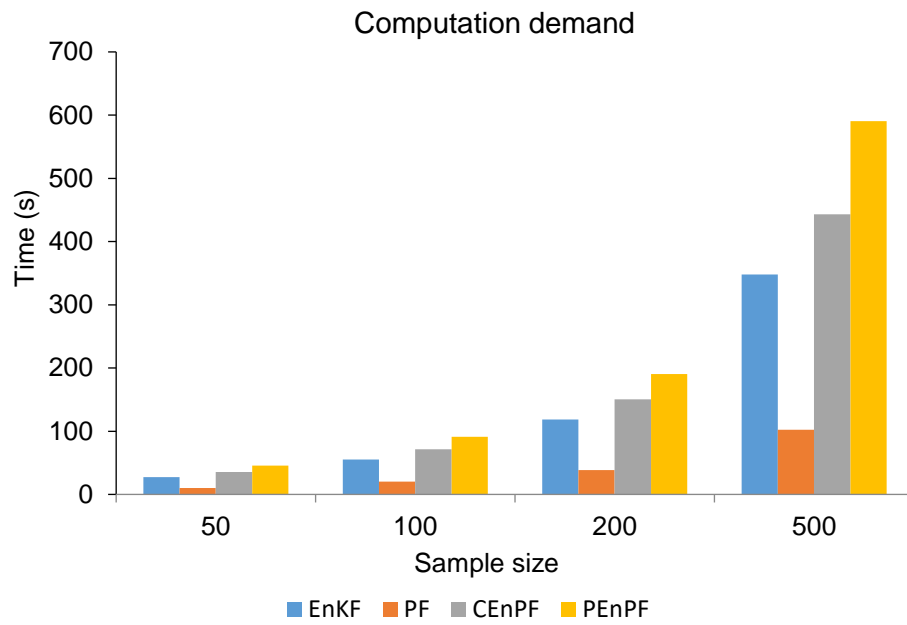


Figure 11. Computation demand for EnKF, PF, CEnPF and PEnPF under different sample size scenarios

Table 1. The predefined true values (used in synthetic experiment), initial fluctuating ranges of Hymod parameters

Description	Parameter	Range	Synthetic true value
Maximum storage capacity of watershed	C_{max} (mm)	[100, 700]	428.18
Spatial variability of soil moisture capacity	b_{exp}	[2, 15]	8.79
Factor distributing flow to the quick-flow tank	α	[0.10, 0.70]	0.28
Residence time of the slow-flow tank	R_s (1/day)	[0.001, 0.20]	0.042
Residence time of the quick-flow tank	R_q (1/day)	[0.2, 0.99]	0.79

1

2 Table 2. the location of rain gauge stations in Huanjiang river basin

Name	Longitude	Latitude
Ganjipan	107.22	37.30
Fanxue	107.58	37.08
Shancheng	107.03	36.95
Gengwan	107.27	36.88
Honglaochi	106.78	36.87
Siheyuan	107.45	36.82
Hongde	107.20	36.77

3

4 Table 3 Locations of National meteorological stations in Jing river basin

Name	Longitude	Latitude
Changwu	107.80	35.20
Xifengzhen	107.63	35.73
Guyuan	106.27	36.00
Huanxian	107.30	36.58
Tongchuan	109.07	35.08

5

6

7

8 Table 4. The NSE coefficient between the ensemble predictions and real observations in
 9 Huanjiang river.

		50	100	200	500
CEnPF	Mean	0.7548	0.7803	0.7736	0.8007
	Min	0.2174	0.6047	0.6620	0.6429
	Max	0.8943	0.9044	0.8464	0.9109
PEnPF	Mean	0.6739	0.7175	0.7294	0.7899
	Min	0.6249	0.6613	0.6563	0.7137
	Max	0.7555	0.7702	0.8471	0.8607
EnKF	Mean	0.6532	0.6907	0.7448	0.7181
	Min	0.3035	0.5223	0.6134	0.6738
	Max	0.8140	0.8056	0.7977	0.7667
PF	Mean	0.6470	0.6458	0.6509	0.6660
	Min	0.4521	0.4721	0.4176	0.4885
	Max	0.7656	0.7318	0.8383	0.7633

11

12 Table 5. The RMSE values between the ensemble predictions and real observations in

13 Huanjiang river.

		50	100	200	500
CEnPF	Mean	9.2789	9.0914	9.3338	8.6391
	Min	6.4205	6.1079	7.7408	5.8972
	Max	17.4726	12.4186	11.4827	11.8033
PEnPF	Mean	11.2552	10.4769	10.1872	9.0089
	Min	9.7672	9.4682	7.7224	7.3720
	Max	12.0960	11.4950	11.5790	10.5680
EnKF	Mean	11.3404	10.8787	9.9398	10.4714
	Min	8.5184	8.7083	8.8827	9.5404
	Max	16.4840	13.6516	12.2815	11.2803
PF	Mean	10.5186	10.5716	10.4382	10.2374
	Min	8.6479	9.2499	7.1836	8.6903
	Max	13.2215	12.9784	13.6322	12.7747

14

15 Table 6. The CRPS values between the ensemble predictions and real observations in
 16 Huanjiang river.

17

		50	100	200	500
CEnPF	Mean	2.7980	2.5831	2.7709	2.5238
	Min	2.3589	1.9576	2.1644	2.1624
	Max	4.1678	3.0720	3.1563	3.0222
PEnPF	Mean	2.4414	2.2300	2.2268	1.9614
	Min	2.0791	1.9265	1.6249	1.6750
	Max	2.6434	2.5651	2.5963	2.1885
EnKF	Mean	3.3559	2.5764	2.3244	2.4289
	Min	2.1443	2.0683	2.2054	2.2345
	Max	5.2723	3.7094	2.7044	2.6382
PF	Mean	3.9765	4.0262	4.1305	4.2854
	Min	2.9877	2.7904	2.5652	3.2007
	Max	5.4238	4.8530	5.0780	5.5043

18

19

20 Table 7. The Resolution between the ensemble predictions and real observations in Huanjiang
 21 river.

22

		50	100	200	500
CEnPF	Mean	52.4690	48.8849	42.4754	43.7232
	Min	43.2976	39.0868	36.1500	38.7363
	Max	66.7200	62.8025	46.6733	57.6743
PEnPF	Mean	19.4104	17.2911	17.6186	16.6493
	Min	17.5940	14.0080	16.0280	15.0580
	Max	20.9610	19.4370	18.6260	18.6290
EnKF	Mean	35.9948	29.0739	24.6598	21.9759
	Min	28.9328	25.4233	23.6961	21.0699
	Max	42.5571	31.6062	25.1039	22.7798
PF	Mean	41.5654	39.6750	39.8738	38.5949
	Min	33.4221	33.5924	21.0764	31.9602
	Max	48.1742	49.7531	55.9405	45.8325

23

Table 8. Comparison of different data assimilation approaches at Xingxi River

		NSE	RMSE	CRPS	Resolution
50	EnKF	0.5553	43.9565	15.2674	23.5072
	PF	0.6837	36.4071	19.0750	32.4610
	CEnPF	0.6951	36.3942	18.4432	39.8297
	PEnPF	0.7294	33.6750	21.2260	24.2767
100	EnKF	0.6014	41.6133	14.1384	21.8007
	PF	0.7338	34.0062	18.5035	23.0801
	CEnPF	0.7127	35.3301	17.1706	24.2102
	PEnPF	0.7166	35.0884	21.0474	12.9162
200	EnKF	0.6110	41.1089	13.8818	20.8912
	PF	0.7163	34.4767	19.5430	19.4740
	CEnPF	0.6725	37.7190	17.6068	21.2002
	PEnPF	0.7465	33.1868	16.8556	21.7079
500	EnKF	0.5231	45.5183	14.8714	22.2468
	PF	0.6786	36.6998	18.6901	22.2949
	CEnPF	0.7530	32.7555	15.8585	20.2561
	PEnPF	0.7403	32.9869	15.7859	24.3501

25

26

27

28

29