

1 **Improving streamflow predictions at ungauged locations**
2 **with real-time updating: Application of an EnKF-based**
3 **state-parameter estimation strategy**

4 **Xianhong Xie¹, Shanshan Meng¹, Shunlin Liang^{1,2}, Yunjun Yao¹**

5 [1] State Key Laboratory of Remote Sensing Science, College of Global Change and Earth
6 System Science, Beijing Normal University, Beijing, China

7 [2] Department of Geographical Sciences, University of Maryland, College Park, USA

8 Correspondence to: Xianhong Xie (xianhong@bnu.edu.cn)

1 **Abstract**

2 The challenge of streamflow predictions at ungauged locations is primarily attributed to
3 various uncertainties in hydrological modelling. Many studies have been devoted to
4 addressing this issue. The similarity regionalization approach, a commonly used strategy, is
5 usually limited by subjective selection of similarity measures. This paper presents an
6 application of a **partitioned** update scheme based on the ensemble Kalman filter (EnKF) to
7 reduce the prediction uncertainties. This scheme performs real-time updating for states and
8 parameters of a distributed hydrological model by assimilating gauged streamflow. The
9 streamflow predictions are constrained by the physical rainfall-runoff processes defined in the
10 distributed hydrological model and by the correlation information transferred from gauged to
11 ungauged basins. This scheme is successfully demonstrated in a nested basin with real-world
12 hydrological data where the subbasins have immediate upstream and downstream neighbours.
13 The results suggest that the assimilated observed data from downstream neighbours have
14 more important roles in reducing the streamflow prediction errors at ungauged locations. The
15 real-time updated model parameters remain stable with reasonable spreads after short-period
16 assimilation, while their estimation trajectories have slow variations, which may be
17 attributable to climate and land surface changes. Although this real-time updating scheme is
18 intended for streamflow predictions in nested basins, it can be a valuable tool in separate
19 basins to improve hydrological predictions by assimilating multi-source datasets, including
20 ground-based and remote-sensing observations.

21

1 **1 Introduction**

2 The streamflow prediction plays a central role in hydrology because it is an important
3 element for water resources management, the design of hydraulic infrastructures and flood
4 risk mapping (Srinivasan et al., 2010). Because it is an important component in the terrestrial
5 water budget, streamflow is also a direct diagnostic variable measuring the impact of climate
6 changes and human activities that act on a given watershed. Streamflow prediction depends
7 highly on reliable hydrological data and sophisticated hydrological models. However,
8 hydrological data are often insufficient due to ungauged or poorly gauged basins in many
9 parts of the world (Sivapalan, 2003). Because of the scarcity of data, hydrological modelling
10 is also plagued by various sources of uncertainties. To reduce uncertainties from those
11 hydrological data and hydrological modelling, the International Association of Hydrological
12 Sciences (IAHS) launched an initiative on Predictions in Ungauged Basins (PUB) (Sivapalan,
13 2003; Sivapalan et al., 2003).

14 Through the past PUB decade, major advances have been achieved including data
15 acquisition and exploitation, modelling strategies and uncertainty analysis, and catchment
16 classification and new theory (Hrachowitz et al., 2013). There is a growing consensus that
17 remote sensing techniques provide valuable data for understanding the land surface
18 hydrological system (Yang et al., 2013). Moreover, considerable progress has been made on
19 hydrological models (typically the distributed hydrological models) to capture the physical
20 process associated with the basin rainfall–runoff and snowmelt–runoff responses. This
21 progress has fostered specific problem areas in the field: uncertainty quantification with
22 respect to model input forcing, model structures and parameters (Ajami et al., 2007; Vrugt et
23 al., 2008; Gupta et al., 2012). To reduce the uncertainty from model parameters, one common
24 practice is the parameter calibration by adjusting model parameters to make the simulated
25 water discharges correspond to the observations (typically the data from the outlet of a
26 watershed) (Duan et al., 1992; Duan et al., 1994). However, a calibrated parameter set with
27 acceptable streamflow simulation performance at the watershed outlet does not guarantee the

1 performance at interior locations (Zhang et al., 2008).

2 The essence of PUB is to transfer information from neighbouring basins to the basins of
3 interest (Sivapalan et al., 2003). Such process is generally referred to as hydrological
4 regionalization, based on either regression methods or measureable distances (with respect to
5 physical similarity or spatial proximity) between gauged and ungauged locations (Hrachowitz
6 et al., 2013). Regionalization techniques regarding model parameters are popular for
7 discharge prediction in ungauged basins. Merz and Blöschl (2004) evaluated the performance
8 of various regionalization methods for parameters of a conceptual catchment model,
9 determining that spatial proximity is able to represent the unknown controls on the runoff
10 regime and the relationships of model parameters within neighbouring basins. Sellami et al.
11 (2013) presented a model parameter regionalization approach based on physical similarity
12 between gauged and ungauged catchments, indicating that similar hydrological behaviour
13 may appear due to physically similar catchments in the same geographic and climatic region.
14 Parajka et al. (2013) reported that the spatial proximity and geostatistics probably perform
15 better than the regression or regionalization with a simple averaging of model parameters
16 from gauged catchments. One drawback of the regionalization of model parameters is that it
17 often confronts an arbitrary criterion for selecting the “behavioural” model parameter sets
18 from the gauged catchment (Sellami et al., 2013). Hrachowitz et al. (2013) provides a
19 comprehensive review of the parameter regionalization and catchment similarity.

20 In addition to those parameter regionalization approaches, newly developed data
21 assimilation methods are also encouraging and are capable to address some issues associated
22 with PUB. They are generally based on physical correlations between the neighbouring basins,
23 and they can combine multi-source observations to transfer information from gauged to
24 ungauged basins (Sivapalan et al., 2003; Troch et al., 2003; Chen et al., 2011). As a typical
25 sequential data assimilation approach, ensemble Kalman filter (EnKF) is popular in hydrology
26 (Reichle et al., 2002; Evensen, 2003; Evensen, 2009). EnKF is attractive in hydrology
27 primarily because it can perform real-time updating with simple implementation and it

1 considers various uncertainties in modelling and observations (Blöschl et al., 2008). The
2 feature of real-time updating is very important for flood forecasting (Norbiato et al., 2008). In
3 some current applications, EnKF is mainly dedicated to dynamic state estimations in which
4 the model parameters are defined with prior values or calibrated in advance (Vrugt et al., 2005;
5 Clark et al., 2008).

6 The EnKF method also provides a general framework to perform state-parameter
7 estimation which is the core of PUB issues. It has been successfully used for parameter
8 estimation of hydrological models. Moradkhani et al. (2005b) proposed a dual state-parameter
9 estimation of hydrological models and made an acceptable application of this method for a
10 lumped hydrological model. Wang et al. (2009) presented three constrained schemes with
11 EnKF to prevent the violation of parameter physical constraints. Most of these studies
12 performed parameter estimations for lumped hydrological models with a small number of
13 parameters to be estimated. Xie and Zhang (2010) successfully demonstrated a joint
14 state-parameter estimation based on EnKF for a distributed hydrological model, i.e., Soil and
15 Water Assessment Tool (SWAT), focusing on one dominant parameter in SWAT. For multiple
16 types of parameter estimation, Xie and Zhang (2013) developed a [partitioned](#) update scheme
17 and indicated the potential of this scheme for streamflow predictions in ungauged basins
18 based on distributed hydrological models.

19 In this study, we present the application of the [partitioned](#) update scheme to improve
20 streamflow predictions in ungauged locations by assimilating gauged streamflow. This data
21 assimilation algorithm is fully coupled with the distributed hydrological model, i.e., SWAT.
22 The state vector and parameters in ungauged subbasins are estimated when information is
23 transferred from gauged subbasins. To our knowledge, this study is the first one which
24 explicitly employs a data assimilation method with state-parameter estimation to improve
25 streamflow predictions in ungauged locations. Although a few applications of data
26 assimilation methods are dedicated to streamflow predictions based on distributed models
27 (Clark et al., 2008; Chen et al., 2011; Lee et al., 2012; Rakovec et al., 2012; McMillan et al.,

1 2013), the model parameter estimation, which is important for PUB, is not systematically
2 considered. In addition to the EnKF-based scheme, note that the other data assimilation
3 methods, e.g., the particle filter (Moradkhani et al., 2005a; DeChant and Moradkhani, 2012),
4 the Particle-DREAM (Vrugt et al., 2013) and the Maximum Likelihood Ensemble Filter (Tran
5 et al., 2014), may also be optional for state-parameter estimation.

6 In the following sections, we first introduce the EnKF-based data assimilation scheme and
7 give a brief description of the SWAT model. We then present an application case concerning a
8 real-world problem in the Zhanghe River basin in China in which river channels are
9 connected and subbasins have nested upstream and downstream neighbours. Three scenarios
10 regarding different combinations of observed streamflow are designed to discuss the impact
11 of gauged locations on streamflow predictions. Finally, conclusions are given in the last
12 section.

13 **2 Methodology**

14 **2.1 EnKF-based state and parameter estimation scheme**

15 To describe the information transfer process from gauged to ungauged locations, we define
16 a joint state vector X that contains gauged (x_g) and ungauged (x_u) states: $X = [x_g, x_u]$.
17 Moreover, we consider the diagnostic variables, i.e., the water discharge and the
18 evapotranspiration, as model states and include them in the vector X to perform streamflow
19 updating in the data assimilation. The joint state vector X and the parameter vector θ
20 estimation at time t are conditioned on measurements (y_t) from gauged basins. The
21 information transfer process, i.e., the posterior probability density function (pdf) $p(X_t, \theta_t | y_t)$,
22 can be expressed within Bayes' framework,

$$23 \quad p(X_t, \theta_t | y_t) \propto p(y_t | X_t, \theta_t) \cdot p(X_t, \theta_t | X_{t-1}, \theta_{t-1}), \quad (1)$$

24 where $p(y_t | X_t, \theta_t)$ is the likelihood function of measurements given model estimations at

1 time t . Moreover, $p(X_t, \theta_t | X_{t-1}, \theta_{t-1})$ is the prior pdf of X and θ at time t that represents
 2 model forecasting and parameter evolutions.

3 The updating framework defined in equation (1) is well included in and effectively solved
 4 by sequential data assimilation strategies, typically, the EnKF strategy (Evensen, 1994). The
 5 EnKF strategy operates sequentially with a forecast step and a filter update step. In the
 6 forecasting process, uncertainty propagation is characterised by an ensemble of model
 7 realisations:

$$8 \quad X_t^{i-} = M(X_{t-1}^{i+}, \theta_{t-1}^{i-}, u_t^i) + \omega_t^i, \quad \omega_t^i \sim N(0, W_t), \quad i = 1, 2, \dots, N, \quad (2)$$

9 where “-” and “+” denote the forecast and analysis for the state vectors X and the parameter
 10 vector θ , t is the time step, u is the input forcing vector, and N is the ensemble size. The model
 11 error vector ω is assumed to follow a Gaussian distribution with zero mean and covariance
 12 W_t . Equation (2) is a general expression with representative errors for all state variables. In
 13 implementation, one may define errors for only a few of the state variables (e.g., soil moisture)
 14 to reflect realistic modeling uncertainties. Detailed prescription of the errors will be given in
 15 section 3.2.

16 Prior to model forecasting using equation (2), the model parameters can be perturbed,
 17 similar to the forecast of the state vector, to avoid the shrinkage of the parameter ensemble
 18 during the updating (Wang et al., 2009). However, the parameter perturbation is susceptible to
 19 over-dispersion in sampling (Moradkhani et al., 2005b). A kernel smoothing technique is
 20 effective to address the over-dispersion while maintaining a reasonable ensemble spread for
 21 the parameters (Liu, 2000; Moradkhani et al., 2005b; Xie and Zhang, 2013). This technique is
 22 briefly expressed as

$$23 \quad \theta_t^{i-} = \alpha \theta_{t-1}^{i+} + (1 - \alpha) \bar{\theta}_{t-1}^+ + \tau_t^i, \quad \tau_t^i \sim N(0, T_t), \quad (3)$$

$$24 \quad \bar{\theta}_{t-1}^+ = \frac{1}{N} \sum_{i=1}^N \theta_{t-1}^{i+}, \quad (4)$$

$$1 \quad T_t = h^2 \text{var}(\theta_{t-1}^+), \quad (5)$$

2 where α is the shrinkage factor typically within [0.95, 0.99], h is the smoothing factor, and T_t
3 is the covariance constrained by the ensemble variance $\text{var}(\theta_t^+)$. The smoothing factor h is
4 defined as $\sqrt{1-\alpha^2}$ to maintain equal variances of the parameter before and after the
5 perturbation. This kernel smoothing technique has been discussed based on synthetic cases
6 (Liu, 2000; Moradkhani et al., 2005b; Xie and Zhang, 2013), so we do not provide any more
7 experiments to demonstrate the properties of the kernel smoothing. The prescription of the
8 shrinkage factor α is subject to trial and error experimentation, but it has limited impact on the
9 parameter estimation (An illustrative case was shown in the response to the reviewers'
10 comments at version 4 of this paper). In this study, it is specified with 0.98 according to the
11 suggestions by Moradkhani et al., (2005b) and Xie and Zhang, (2013).

12 With the forecast of the states and parameters, the filter update step is performed when
13 observations are available. This updating is actually the solving process for equation (1). Here
14 we intentionally create an explicit expression of the updating for gauged and ungauged states
15 and parameters:

$$16 \quad \begin{bmatrix} x_{g,t}^{i+} \\ x_{u,t}^{i+} \\ \theta_t^{i+} \end{bmatrix} = \begin{bmatrix} x_{g,t}^{i-} \\ x_{u,t}^{i-} \\ \theta_t^{i-} \end{bmatrix} + K_t \cdot (y_t^i - Hx_{g,t}^{i+}), \quad (6)$$

17 where y_t^i is the observation vector, which is appropriately perturbed using covariance of R
18 to account for uncertainties in observations, and H is the observation operator and it is linear
19 in this study. The Kalman gain matrix K_t is expressed as

$$20 \quad K_t = \begin{bmatrix} \text{cov}(x_{g,t}, x_{g,t}) \\ \text{cov}(x_{g,t}, x_{u,t}) \\ \text{cov}(x_{g,t}, \theta_t) \end{bmatrix} \cdot (\text{cov}(x_{g,t}, x_{g,t}) + R)^{-1}, \quad (7)$$

21 where $\text{cov}(\cdot)$ is the covariance operator that is computed from the ensembles of states and

1 parameters. Please note the size of the matrix K_t is $n \times m$, where n is the total number of state
2 variables and parameters and m is the number of observations.

3 The above two equations rely on EnKF with a state-augmentation technique. This
4 technique is valid and able to retrieve correct parameter estimates in real time primarily
5 because it allows for parameter dynamics and performs the parameter evolution. Specifically,
6 model parameters are assumed as an extension of state variables and they can travel slowly
7 with time, in response to changes in environmental forcing inputs (Liu and Gupta, 2007). Like
8 the model state forecasting, the parameters are perturbed/evolved using the kernel smoothing
9 technique. In this way, the evolution of model parameters is consistent with the forecasting of
10 model state variables. Thus the model parameters can be appended to the state vector
11 (Moradkhani et al., 2005; Xie and Zhang, 2010, 2013). When observations are available, the
12 parameters are updated along with state variables by assimilating these observations.
13 Therefore, their estimates are expected to converge to the “correct” posterior target
14 distribution (Xie and Zhang, 2013). This technique has been successfully used in many cases
15 for real-time state and parameter estimation (Moradkhani et al., 2005b; Wang et al., 2009; Xie
16 and Zhang, 2010, 2013).

17 Moreover, we can see that EnKF provides a general framework to transfer information
18 from gauged to ungauged basins. However, when used for parameter estimations in
19 distributed hydrological models, it is vulnerable to corruption due to spurious covariance
20 computation in equation (7), primarily resulting from a large degree of freedom for
21 high-dimensional vectors of the augmented state. To relieve this problem, Xie and Zhang
22 (2013) proposed a partitioned forecast-update scheme (PU_EnKF) that is inspired by the dual
23 state-parameter estimation algorithm (Moradkhani et al., 2005b). In the partitioned
24 forecast-update scheme, the parameter set of a hydrological model is partitioned into different
25 types (N_p types in total) based on their sensitivities. Each type is estimated in an individual
26 loop by repeated forecasting and updating. Here, the parameter type maintains an aggregation
27 connotation. A parameter type can contain only one parameter (e.g., for lumped hydrological

1 models) or many parameters associated with the same number of computational units in
2 distributed hydrological models. For example, the parameter CN_2 in SWAT (will be
3 introduced in subsection 2.2) is considered as a parameter type.

4 At time t , the PU_EnKF is iteratively applied as follows for N_p loops:

5 (I) Perform parameter evolution using equation (3) for the j th parameter type, producing a
6 new ensemble of parameters.

7 (II) Run the model N times following equation (2) to obtain ensemble predictions for gauged
8 and ungauged state variables. In the prediction, the j th parameter type is prescribed with a
9 member of the ensemble produced in step (I), while the others are set with the ensemble
10 means that are estimated from previous loops at this time step and from the previous time
11 step.

12 (III) Compute the Kalman gain matrix using equation (7) based on the ensembles of states and
13 parameters when observations become available at time t .

14 (IV) Update the state vector and the j th parameter type using equation (6).

15 (V) Compute the ensemble means of the j th parameter type. The means are the estimates of
16 the parameters and will be used in step (II) in the subsequent loops to estimate the other
17 parameter types.

18 (VI) Return to step (I) if $j < N_p$. Otherwise, go to the next time step $t + 1$. The updated state
19 vector from the loop $j = N_p$ is considered as estimates of gauged and ungauged state variables;
20 and all estimates of parameters are also obtained.

21 We can see that the partitioned update scheme employs an iterative algorithm to update
22 each parameter type at each time step, not only is one parameter considered at a time. At time
23 t , the new estimated parameter values from previous loops are used for the model forecasting
24 (Eq. (2)) in the current loop in which a target parameter type (the j th parameter type) is
25 estimated. This iterative update is expected to push the estimates towards their optimal values.

1 Therefore, this scheme is quite suitable for distributed hydrological models to estimate
2 high-dimensional parameters. Its capability has been demonstrated using synthetic cases and
3 it has been successfully used in a real watershed for state and parameter estimation (Xie and
4 Zhang, 2013). In this study, we apply this scheme to improve the streamflow prediction in
5 ungauged sites and to estimate model parameters.

6 **2.2 Model description**

7 The distributed hydrologic model, SWAT, is a basin-scale hydrological model developed by
8 the USDA Agricultural Research Service (Arnold et al., 1998; Arnold and Fohrer, 2005). In
9 implementation of SWAT, a basin is partitioned into multiple subbasins that are then divided
10 into hydrologic response units (HRUs), which consist of unique land cover, management, and
11 soil characteristics (Neitsch et al., 2001; Gassman et al., 2007). The HRUs are the basic
12 computational units in which the overall hydrologic balance is simulated, including
13 precipitation partitioning, surface runoff generation, evapotranspiration (ET), soil water and
14 groundwater movement.

15 The surface runoff generation is commonly simulated using the Soil Conservation Service
16 (SCS) model (Rallison and Miller, 1981; Ponce et al., 1996). This model has only one
17 parameter, i.e., the curve number at moisture condition II (CN₂), which is also the dominant
18 parameter in SWAT. Actual ET is formulated based on potential ET to account for evaporation
19 from the plant canopy, transpiration, sublimation and evaporation from the soil. The soil water
20 movement is characterised by a storage routing technique that uses the field capacity to
21 dominate redistribution of water between layers. By infiltration or percolation, a fraction of
22 water below the soil profile enters groundwater storage as recharge and is partitioned between
23 shallow and deep aquifers. Base flow from the shallow aquifer is also routed to river channels.
24 Details regarding these processes can be found in the SWAT user's manual (Neitsch et al.,
25 2001).

26 SWAT contains a large number of spatially varying parameter types to be prescribed before
27 hydrologic simulation and prediction. These parameters consist of the surface roughness, soil

1 properties, land-cover pattern and hydraulic conditions of the river channel. Although their
2 default values can be prescribed according to lookup tables, the optimal values must be
3 calibrated on the basis of modelling behaviour and observations. To reduce the number of
4 calibrating parameters, a sensitivity analysis is usually required (van Griensven et al., 2006).
5 Considerable effort has been devoted to sensitivity analysis for SWAT; several parameters are
6 recognised as the most influential ones that dominate the model behaviour (Holvoet et al.,
7 2005; Muleta and Nicklow, 2005; van Griensven et al., 2006). Based on these studies, seven
8 parameters (also called parameter types) are selected and shown in Table 1. They underpin
9 different hydrologic processes in a basin involving the surface runoff, soil water, baseflow,
10 groundwater, evapotranspiration and channel water processes. Their ranges are determined in
11 terms of the lookup tables (Neitsch et al., 2001) and the specific soil and land use properties
12 of the Zhanghe River basin (Post and Jakeman, 1999).

13 In addition to these sensitive parameter types, ten hydrologic variables are selected to be
14 updated in data assimilation (Table 2). They can be divided into three groups: (1) Quick water
15 storage (marked with QW in Table 2) regarding surface runoff, (2) Slow water storage
16 (marked with SW) associated with baseflow and groundwater flow and soil moisture, and (3)
17 river channel storage (marked with CW) and flow. The first nine variables are the dynamic
18 states that characterise water storage status in HRUs or subbasins and partially influence the
19 diagnostic variables, i.e., ET and the water discharge (Q_r). Therefore, along with both outputs,
20 these states should be updated to guarantee consistent model behaviour. In this study, ET is
21 excluded from the state vector because there are no ET observations and its passive update in
22 data assimilation does not impact other state estimations.

23 The SWAT model is used for this study for two main reasons. First, SWAT is a very popular
24 distributed hydrological model to predict water, sediment, and agricultural chemical yields in
25 large, complex watersheds (Gassman et al., 2007). An improved version of this model has
26 been used to simulate the water movement in the Zhanghe River basin, an irrigation district
27 with paddy rice planting (Xie and Cui, 2011). Second, we have coupled it with the

1 EnKF-based algorithms with a few successful applications (Xie and Zhang, 2010; Xie, 2013;
2 Xie and Zhang, 2013). Therefore, such a coupled SWAT-EnKF data assimilation platform is
3 expected to be more powerful and widely used for real-time hydrological predictions. SWAT
4 requires a significant amount of data including model input and system response data (e.g.,
5 streamflow, evapotranspiration), which seems not consistent with effort of predictions in
6 ungauged basins. But this issue can be eased to some degree because streamflow data from
7 just a few locations at downstreams (e.g. the outlet) can favour estimation for the entire basin
8 by the data assimilation scheme used in this study.

9 **3 Application to a real case**

10 **3.1 Study area and database**

11 The data assimilation scheme is applied in the Zhanghe River basin in Hubei Province,
12 China (Figure 1). The Zhanghe drains an area of 1129 km², and the elevation difference
13 between the north and the south is more than 400 m. It has a typical subtropical climate with
14 an annual mean temperature of 17 °C. The annual rainfall in the catchment is approximately
15 970 mm per year, although rainfall varies substantially from year to year depending upon the
16 monsoon strength. This basin is actually an agricultural irrigation area and its cultivated area
17 accounts for 59%. Paddy rice is the primary cultivated plant, which, from May to August,
18 requires irrigation water from the Zhanghe reservoir and thousands of local ponds. Owing to
19 intense human activities, including cultivation, irrigation and drainage, streamflow prediction
20 in this basin is challenge with large uncertainties (Cai, 2007; Xie and Cui, 2011).

21 We choose the Zhanghe River basin as a study area because there are relatively sufficient
22 datasets associated with weather conditions, land use and soil properties, and hydrological
23 information. This area has been chosen for a few modelling studies (Cai, 2007; Xie and Cui,
24 2011). The land use classification with resolution of 14.25 m was retrieved based on remote
25 sensing data (Landsat ETM+) for years 2000 and 2001 (Figure 1 (b)). The land use pattern in
26 this basin exhibits only small changes since 2000. Therefore, we assume the land use pattern

1 in the period 2004-2006 is the same as in 2000-2001. The soil map with soil properties, which
2 is used to derive model parameters, is obtained from the local agriculture department. The
3 weather dataset, including daily temperature, radiation, wind speed and relative humidity,
4 from January 2000 to December 2006 is available from five stations distributed in and around
5 this basin as shown in Figure 1 (c). Moreover, four streamflow gauges were installed, marked
6 as A, B, C and D for simple referencing. Gauge D is the outlet of the basin. Gauge A is
7 located at the outlet of a small source subbasin. Because these four gauges observe the river
8 stages and then transform the data into streamflow according to calibrated rating curves, daily
9 streamflow data for the period 2003-2006 are available.

10 The Zhanghe River basin is divided into 20 subbasins based on a digital elevation model
11 (DEM) with a resolution of 90 m (Figure 1 (c)). Thereafter, 98 HRUs are obtained according
12 to land use and the soil map. With this delineation, Gauge A drains runoff from a source
13 subbasin, Gauge B drains four, Gauge C drains ten, and Gauge D drains all the basins.

14 **3.2 Error quantification**

15 The success of ensemble-based data assimilation methods depends partly on ensemble
16 generations to quantify errors from model input forcing, parameters and model structures.
17 Moreover, quantifying observation errors is also critical to account for uncertainties from
18 measurements and derivations. Due to the dynamics of the SWAT model, the
19 errors/uncertainties from the input forcing, parameters and the model structure are transferred
20 to the water storages (e.g., soil moisture and channel storages) and diagnostic variables (e.g.,
21 streamflow). Although ten selected variables require updating in SWAT, two of them are
22 perturbed in this study to represent the modelling uncertainties, i.e., soil moisture and
23 streamflow, because the other variables are internal and their uncertainties are transferred to
24 the soil moisture and the simulated streamflow (Xie and Zhang, 2013). Moreover,
25 precipitation as a major forcing input is also perturbed to represent the uncertainty probably
26 derived from weather forecasting and other sources.

27 Perturbations to the above three variables are conducted based on zero-mean Gaussian

1 distributions. The standard deviation (σ) for SWAT-simulated soil moisture is set as 0.03
2 m^3/m^3 as suggested by Chen et al. (2011). The standard deviations for streamflow and
3 precipitation are assumed to be proportional to their values (Clark et al., 2008),

$$4 \quad \sigma_x = \eta_x \cdot x, \quad (8)$$

5 where η is the fractional factor of the standard deviation to the variable x . Thus, there are three
6 fractional factors corresponding to the simulated streamflow (η_{Q_m}), observed streamflow (η_{Q_o})
7 and precipitation (η_p). Therefore the PU_EnKF scheme used in this study is also applicable to
8 hydrological prediction when measured rainfall data is unavailable but could be derived from
9 various sources (e.g., weather forecasting). With this error quantification, the three standard
10 deviations vary with time, depending on the magnitudes of the four variables.

11 These fractional factors should not only represent the related uncertainties in modelling and
12 the observations but also produce ensemble streamflow predictions with reasonable ensemble
13 spread (Clark et al., 2008). Based on the uncertainty analysis by Xie and Cui (2011), the
14 prediction errors with the SWAT model are more than 10% of the variables due to the
15 irrigation and drainage practices in the Zhanghe River basin; the measurement of precipitation
16 also has the same level of uncertainty. Therefore, various combinations of factor values are
17 evaluated by running the data assimilation procedure. Table 3 presents the final choice of the
18 **three** fractional factors.

19 Note the error quantification remains challenging for land surface data assimilation. A few
20 newly developed approaches may be a good attempt, e.g., adaptive filtering (Crow and
21 Reichle, 2008; Reichle et al., 2008). However, we quantify the model and observation
22 uncertainties in terms of an experiential and practical perspective in which large storm events
23 normally induce larger uncertainties in modelling and observations. Moreover, an
24 overestimation of uncertainties is a better practice than underestimation to avoid the ensemble
25 shrinkage (Crow and Van Loon, 2006; Clark et al., 2008).

1 3.3 Assimilation setup and scenario design

2 The assimilation process is performed with three successive periods (Xie and Zhang, 2013).
3 First, the model is prescribed with prior parameters and spun-up within the period 1/1/2003 to
4 6/30/2003 to initialise the model states. At the end of this period, the seven parameters of the
5 SWAT model are perturbed using the Latin hypercube method (Helton and Davis, 2003) with
6 Gaussian distributions. The parameter means [regarding](#) the Gaussian distributions are set
7 according the lookup table suggested in SWAT (Neitsch et al., 2001); the associated variances
8 are constrained to ensure that random samples are within their respective physically or
9 model-required ranges in Table 2. Please note the uniform distribution is more intuitive than
10 the Gaussian and often also used in sampling (Moradkhani et al., 2005b). In this study, we use
11 the Gaussian because the lookup table provides prior estimates for the parameters. The
12 number of parameter samples (i.e., the ensemble size) is 80. After the parameter perturbations,
13 the second period begins (7/1/2003 – 12/31/2003) to perturb the model input forcing, model
14 states and diagnostic variables as described in subsection 3.2. The aim of this perturbation
15 period is to quantify the uncertainties in prediction and to generate reasonable ensemble
16 spread for subsequent data assimilation. The third period is the data assimilation period
17 (1/1/2004 – 12/31/2005) in which the streamflow observations are assimilated when data are
18 available. Given that streamflow originates primarily from either surface runoff or subsurface
19 runoff in different periods, the variables of quick water storage (QW in Table 2) are updated
20 only when precipitation occurs. The variables of slow water storage (SW) are updated during
21 dry periods (no precipitation), and variables of channel water storage (CW) are updated at
22 every time step.

23 To demonstrate the improvement of streamflow prediction in ungauged locations, we only
24 assimilate streamflow from one or two of the four gauges and the remaining gauges, regarded
25 as pseudo-ungauged locations, are used to validate the performance of data assimilation.
26 Three scenarios with different combinations of data from the four gauges are designed:

27 (I) ASS_D: The observed data of streamflow from Gauge D are assimilated; Gauges A, B

1 and C are assumed as pseudo-ungauged. This scenario is similar to a common calibration
2 practice for which only the outlet (Gauge D) discharge data are employed to calibrate the
3 parameters and to extrapolate streamflow of ungauged subbasins.

4 (II) ASS_{BD}: The observed data of streamflow from Gauge B and D are assimilated; the
5 other two are regarded as pseudo-ungauged subbasins. This scenario adds the data from
6 Gauge B at the upstream in this basin based on scenario ASS_D.

7 (III) ASS_{AB}: The observed data of streamflow only from Gauge A and B are assimilated.
8 This scenario only uses the streamflow from the two gauges in the upstream subbasins.

9 **3.4 Prediction in ungauged locations**

10 Ensemble streamflow predictions along with parameter estimations are performed for the
11 three scenarios. To distinguish the improvement of streamflow prediction, a control-run
12 scenario is conducted in which the model parameters are prescribed with the calibrated
13 estimates from Xie and Cui (2011). The data assimilation performance is evaluated by
14 comparing with the four series of observed streamflow. Although the observed streamflow
15 series still contain uncertainties, we consider them to be a benchmark because the
16 observations are commonly assumed to be the best estimates of “real” streamflow processes.
17 Therefore, the series of streamflow prediction errors are computed (predictions minus
18 observations). The root-mean-square error (RMSE) and the mean absolute error (MAE) are
19 used as comprehensive indexes for evaluations. To quantify the ensemble spread of
20 streamflow in data assimilation, we define a measure, i.e., ensemble coverage index (EnCI)
21 that is a percent of discharge data contained in the 95% ensemble simulation intervals.

22 Figure 2 shows the streamflow errors from the control-run prediction and scenario ASS_D.
23 The reason the errors being presented instead of the streamflow observations is that some of
24 the streamflow observations are so large that the difference between the cases is not notable.
25 The control-run simulation clearly overestimates the peak flow (in wet periods of rainfall
26 occurrence) for the four gauges, while underestimates the base flow in some dry periods (e.g.,
27 230th–300th time steps). This poor performance is significantly improved by assimilating the

1 observed streamflow and by considering the uncertainties from the input forcing and model
2 states. It may not be surprising that the Gauge D streamflow errors in ASS_D are less than
3 those in the control-run scenario because the observed streamflow from Gauge D is
4 assimilated to update the prediction. For the (pseudo-) ungauged locations, the streamflow
5 predictions of Gauge A, B and C are also more acceptable than from the control-run scenario.
6 At Gauge C, for example, the RMSE decreases from 3.539 m³/s to 1.912 m³/s. Moreover,
7 there is no notable biased prediction due to the slight overestimations and underestimations
8 for peak flow.

9 The EnCI for Gauge D is up to 95.72% (see Figure 2). This means that 95.72% discharge
10 data are contained in the 95% ensemble intervals, except that some discharge data with
11 considerable magnitudes of flood are outside of the intervals. The lowest EnCI for Gauge A
12 (75.21%) is partly due to the fact that Gauge A is the farthest gauge to the outlet (Gauge D, its
13 data are assimilated). Nevertheless, all ensemble spreads for the four gauges are reasonable to
14 trace and to contain the discharge data.

15 Figure 3 shows the results for Gauge C from scenarios ASS_BD and ASS_AB. Adding an
16 observed gauge (Gauge B) at the upstream in the basin, i.e., the ASS_BD scenario, provides
17 better streamflow predictions in the pseudo-ungauged subbasins than the ASS_D scenario; the
18 RMSE drops to 1.669 m³/s and the EnCI is up to 90.28%. If assimilating the data from the
19 upstream locations, i.e., the ASS_AB scenario, the improvement is degraded and the
20 predictions are only slightly better than the control-run scenario. The improvement of
21 streamflow prediction using the PU_EnKF scheme depends on the correlation of physical
22 processes between gauged and ungauged locations. If the two locations are very close (which
23 means the correlation of flow processes will be strong), quit favorable data assimilation
24 performance will be shown. In addition to Gauge C (for pseudo-ungauged locations), Gauge
25 A, B and D have encouraging streamflow predictions due to the fact the data from these
26 gauges are assimilated to update the predicted streamflow (not shown in Figure 3).

27 Along with the updating of model states and diagnostic variables, the model parameters are

1 also estimated. Figure 4 shows examples of real-time parameter updating from the ASS_D
2 scenario. After about 130 time steps, the ensemble trajectories are nearly stable with slow
3 variations which are probably induced by the changes of land surface and river channel
4 conditions for runoff generation and routing (Liu et al., 2008; Troch et al., 2013). At every
5 time step in data assimilation, the parameter samples can be approximated with Gaussian
6 distributions and they are constrained within the prior ranges (Min – Max, see Table 1) as
7 shown in the histograms in Figure 4. This property is favourable for parameter estimation
8 with ensemble-based data assimilation. The uncertainties of parameter estimates at every time
9 step are represented using the ensemble spread (EnSp), which is computed based on sample
10 variances (see the illustration under Figure 5). At the beginning of the data assimilation, the
11 parameters have broad ensemble spreads. The spreads quickly shrink after 100 time steps with
12 the evolution of the streamflow assimilation, and remain stable after 400 time steps. Therefore,
13 the estimate uncertainties of the parameters decrease with the data assimilation and state
14 updating. Moreover, the relative stabilities of ensemble trajectories (Figure 4) and the
15 ensemble spreads (Figure 5) imply an attractive potential that it is possible to use short-term
16 data to retrieve optimal estimates of parameters.

17 Even though the three scenarios provide different parameter estimates due to the
18 assimilation of different observations, encouraging properties of parameter estimations are
19 achieved in the three scenarios. It is not sure so far whether the parameter estimates converge
20 to their appropriate values in this real-world application, so the parameter estimates require a
21 further validation to evaluate the effectiveness of the PU_EnKF scheme.

22 **3.5 Validation for parameter estimates**

23 It is difficult to directly validate the parameter estimates using measurements because the
24 SWAT model is a conceptual hydrological model and most parameters do not have physical
25 meanings. Only a few parameters (e.g., the SOL_AWC in Table 1) can be measured at local
26 sites; those parameters regarding HRUs, subbasins and river channels remain difficult to be
27 obtained by sampling experiments. We perform single-run predictions using the parameter

1 estimates from the three scenarios and evaluate the predicted streamflow against observed
2 streamflow. This is a commonly used strategy to validate parameters of a conceptual
3 hydrological model. For simplicity and consistency, the three single-run predictions are
4 named ASS_D, ASS_BD and ASS_AB, although they are neither assimilation-based
5 predictions nor ensemble predictions. Moreover, the control-run prediction is used for
6 comparison. All four scenarios are run for the period 1/1/2006 – 10/31/2006. The uncertainties
7 in the input forcing and the model structure are not considered in these predictions.

8 Figure 6 shows the streamflow prediction errors from the four scenarios. Only the results of
9 Gauge C and Gauge D are shown because they are located at the downstream locations in the
10 Zhanghe River basin. The three scenarios using prescribed parameters with estimates from
11 data assimilation achieve better predictions for the two gauges than the control-run scenario.
12 The RMSE of Gauge D from the ASS_D scenario decreases from 5.550 m³/s to 2.324 m³/s.
13 Moreover, the ASS_BD scenario provides the best predictions among the four scenarios. All
14 of these improvements are attributable to the appropriate parameter estimates from the data
15 assimilation. The ASS_BD scenario renders the most reasonable parameter estimates.
16 Comparably, the parameter estimates from ASS_D are also satisfactory for streamflow
17 predictions, while the estimates from the ASS_AB scenario lead to slight improvements for
18 streamflow predictions. Therefore, the parameter estimation performance of the three
19 scenarios is consistent with the prediction of diagnostic variables (i.e., the water discharge) as
20 illustrated in subsection 3.4. The assimilated observations from downstream, especially the
21 outlet of the basin, have more important roles than those from upstream for parameter
22 estimation and streamflow predictions in ungauged subbasins.

23 **4 Conclusions**

24 We present an application of PU_EnKF for improving streamflow predictions at ungauged
25 locations. This scheme features real-time updating and simultaneous state-parameter
26 estimation, considering modelling and observing uncertainties. Moreover, the scheme
27 constrains the predictions by the physical rainfall-runoff processes that are defined in the

1 distributed hydrological model (i.e., the SWAT model), and it accounts for the correlations of
2 states and parameters between gauged and ungauged subbasins. The correlations are
3 represented by the covariance matrix in the Kalman gain. With the constraint and the
4 correlation representation, the observed information is successfully transferred to ungauged
5 locations and thereby improves streamflow prediction.

6 The real-world application case suggests that the PU_EnKF scheme performs better than the
7 control-run simulation (with calibrated parameters) for streamflow predictions at gauged and
8 ungauged locations. Although only the outlet-gauged data are assimilated, the streamflow
9 predictions at ungauged sites are still acceptable, since they contain convergent flow
10 information from all subbasins due to runoff routing. Generally, the downstream data
11 (especially the data from the outlet) have important roles to reflect the runoff generation for
12 the entire basin. This data assimilation scheme provides reasonable estimates of model
13 parameters for all computational units (i.e., subbasins and HRUs), including both gauged and
14 ungauged sites, as validated by the conventional single-run simulation. Moreover, the
15 parameter estimates approach nearly stable levels after a small number of time steps (130
16 steps in this study). The parameter estimates show slow variations that would be an advantage
17 of PU_EnKF to identify the changes of land surface properties.

18 Although favourable performance to improve streamflow predictions is obtained using the
19 EnKF-based scheme, the runoff routing is neglected within the PU_EnKF assimilation setup
20 because the travel time of generated runoff is less than one day in the Zhanghe River
21 watershed. In fact, the time lag of runoff routing is an important factor for short-time (e.g., the
22 hourly step) flood forecasting (Li et al., 2013; Pan and Wood, 2013). Moreover, this scheme is
23 intent on PUB for the nested basins in which the correlations of states and parameters
24 between neighbouring subbasins can be constructed. For separate basins in the same climatic
25 regions and land surface conditions, assimilating other sources of data (e.g., the remotely
26 sensed soil moisture and bright temperature) is expected to improve the predictions of
27 hydrological variables (Troch et al., 2003). Nevertheless, this study provides an encouraging

1 application for PUB by assimilating streamflow, which is generally regarded as quality
2 observations compared with the remote sensing data. There are optional methods to address
3 PUB, e.g., the Particle-DREAM by Vrugt et al., (2013). It will be an encouraging attempt to
4 compare these methods with distributed hydrological models for hydrological diagnosis and
5 predictions.

6 **Acknowledgements**

7 We would like to thank Prof. Giuliano Di Baldassarre and three anonymous reviewers for
8 their constructive comments to polish this paper. Dr. Jasper A. Vrugt provided useful
9 suggestions to improve this study. This work was supported by grants from the National
10 Natural Science Foundation of China (41471019, 51009001), the National High Technology
11 Research and Development Program of China (2013AA121200), and the International S&T
12 Cooperation Program of China (2012DFG21710).

13

1 **References**

- 2 Ajami, N. K., Duan, Q., and Sorooshian, S.: An integrated hydrologic Bayesian multimodel
3 combination framework: Confronting input, parameter, and model structural uncertainty in
4 hydrologic prediction, *Water Resour. Res.*, 43, W01403, doi:10.1029/2005wr004745, 2007.
- 5 Arnold, J. G., and Fohrer, N.: SWAT2000: current capabilities and research opportunities in
6 applied watershed modelling, *Hydrol. Process.*, 19, 563-572, 2005.
- 7 Arnold, J. G., Srinivasan, R., Muttiah, R. S., and Williams, J.: Large area hydrologic modeling
8 and assessment part I: Model development1, *JAWRA J. Am. Water Resour. As.*, 34, 73-89,
9 1998.
- 10 Blöschl, G., Reszler, C., and Komma, J.: A spatially distributed flash flood forecasting model,
11 *Environ. Model. Softw.*, 23, 464-478, doi:10.1016/j.envsoft.2007.06.010, 2008.
- 12 Cai, X. L.: Strategy analysis on integrated irrigation water management with RS/GIS and
13 hydrological model, Ph.D thesis, Wuhan University (China), 2007.
- 14 Chen, F., Crow, W. T., Starks, P. J., and Moriasi, D. N.: Improving hydrologic predictions of a
15 catchment model via assimilation of surface soil moisture, *Adv. Water Resour.*, 34, 526-536,
16 doi:10.1016/j.advwatres.2011.01.011, 2011.
- 17 Clark, M. P., Rupp, D. E., Woods, R. A., Zheng, X., Ibbitt, R. P., Slater, A. G., Schmidt, J., and
18 Uddstrom, M. J.: Hydrological data assimilation with the ensemble Kalman filter: Use of
19 streamflow observations to update states in a distributed hydrological model, *Adv. Water*
20 *Resour.*, 31, 1309-1324, doi:10.1016/j.advwatres.2008.06.005, 2008.
- 21 Crow, W. T., and Van Loon, E.: Impact of incorrect model error assumptions on the sequential
22 assimilation of remotely sensed surface soil moisture, *J. Hydrometeorol.*, 7, 421-432, 2006.
- 23 Crow, W. T., and Reichle, R. H.: Comparison of adaptive filtering techniques for land surface
24 data assimilation, *Water Resour. Res.*, 44, W08423, doi:10.1029/2008wr006883, 2008.
- 25 DeChant, C. M., and Moradkhani, H.: Examining the effectiveness and robustness of

1 sequential data assimilation methods for quantification of uncertainty in hydrologic
2 forecasting, *Water Resour. Res.*, 48, W04518, doi:10.1029/2011wr011011, 2012.

3 Duan, Q. Y., Sorooshian, S., and Gupta, V.: Effective and efficient global optimization for
4 conceptual rainfall-runoff models, *Water Resour. Res.*, 28, 1015-1031, 1992.

5 Duan, Q. Y., Sorooshian, S., and Gupta, V. K.: Optimal Use of the Sce-Ua Global
6 Optimization Method for Calibrating Watershed Models, *J. Hydrol.*, 158, 265-284, 1994.

7 Evensen, G.: Sequential data assimilation with a nonlinear quasi-geostrophic model using
8 Monte Carlo methods to forecast error statistics, *J. Geophys. Res.*, 99, 10143-10162,
9 doi:10.1029/94jc00572, 1994.

10 Evensen, G.: The ensemble Kalman filter: Theoretical formulation and practical
11 implementation, *Ocean. Dynam.*, 53, 343-367, 2003.

12 Evensen, G.: *Data Assimilation: the Ensemble Kalman Filter*, Springer Verlag, Berlin,
13 Heidelberg, 2009.

14 Gassman, P., Reyes, M., Green, C., and Arnold, J.: The soil and water assessment tool:
15 historical development, applications, and future research directions, *T. ASABE*, 50,
16 1211-1250, 2007.

17 Gupta, H. V., Clark, M. P., Vrugt, J. A., Abramowitz, G., and Ye, M.: Towards a
18 comprehensive assessment of model structural adequacy, *Water Resour. Res.*, 48, W08301,
19 doi:10.1029/2011wr011044, 2012.

20 Helton, J., and Davis, F.: Latin hypercube sampling and the propagation of uncertainty in
21 analyses of complex systems, *Reliabil. Eng. Syst. Safe.*, 81, 23-69, 2003.

22 Holvoet, K., van Griensven, A., Seuntjens, P., and Vanrolleghem, P. A.: Sensitivity analysis
23 for hydrology and pesticide supply towards the river in SWAT, *Phys. Chem. Earth, Parts*
24 *A/B/C*, 30, 518-526, doi:10.1016/j.pce.2005.07.006, 2005.

25 Hrachowitz, M., Savenije, H. H. G., Blöschl, G., McDonnell, J. J., Sivapalan, M., Pomeroy, J.

1 W., Arheimer, B., Blume, T., Clark, M. P., Ehret, U., Fenicia, F., Freer, J. E., Gelfan, A., Gupta,
2 H. V., Hughes, D. A., Hut, R. W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P. A.,
3 Uhlenbrook, S., Wagener, T., Winsemius, H. C., Woods, R. A., Zehe, E., and Cudennec, C.: A
4 decade of Predictions in Ungauged Basins (PUB) – a review, *Hydrolog. Sci. J.*, 58, 1198-1255,
5 doi:10.1080/02626667.2013.803183, 2013.

6 Lee, H., Seo, D.-J., Liu, Y., Koren, V., McKee, P., Corby, R., and Pappenberger, F.: Variational
7 assimilation of streamflow into operational distributed hydrologic models: effect of
8 spatiotemporal scale of adjustment, *Hydrol. Earth Syst. Sci.*, 16, 2233–2251, 2012

9 Li, Y., Ryu, D., Western, A. W., and Wang, Q. J.: Assimilation of stream discharge for flood
10 forecasting: The benefits of accounting for routing time lags, *Water Resour. Res.*, 49,
11 1887-1900, doi:10.1002/wrcr.20169, 2013.

12 Liu, F.: Bayesian time series: analysis methods using simulation-based computation Ph.D
13 thesis, Institutes of Statistics and Decision Science, Duke University, Durham, North Carolina,
14 USA, 2000.

15 Liu, G., Chen, Y., and Zhang, D.: Investigation of flow and transport processes at the MADE
16 site using ensemble Kalman filter, *Adv. Water Resour.*, 31, 975-986, 2008.

17 McMillan, H. K., Hreinsson, E. Ö., Clark, M. P., Singh, S. K., Zammit, C., and Uddstrom, M.
18 J.: Operational hydrological data assimilation with the recursive ensemble Kalman filter,
19 *Hydrol. Earth Syst. Sci.*, 17, 21-38, doi:10.5194/hess-17-21-2013, 2013.

20 Merz, R., and Blöschl, G.: Regionalisation of catchment model parameters, *J. Hydrol.*, 287,
21 95-123, doi:10.1016/j.jhydrol.2003.09.028, 2004.

22 Moradkhani, H., Hsu, K.-L., Gupta, H., and Sorooshian, S.: Uncertainty assessment of
23 hydrologic model states and parameters: Sequential data assimilation using the particle filter,
24 *Water Resour. Res.*, 41, W05012, doi:10.1029/2004wr003604, 2005a.

25 Moradkhani, H., Sorooshian, S., Gupta, H. V., and Houser, P. R.: Dual state-parameter
26 estimation of hydrological models using ensemble Kalman filter, *Adv. Water Resour.*, 28,

1 135-147, 2005b.

2 Muleta, M. K., and Nicklow, J. W.: Sensitivity and uncertainty analysis coupled with
3 automatic calibration for a distributed watershed model, *J. Hydrol.*, 306, 127-145,
4 doi:10.1016/j.jhydrol.2004.09.005, 2005.

5 Neitsch, S., Arnold, J., Kiniry, J., Williams, J., and King, K.: Soil and water assessment tool
6 theoretical documentation version 2000, Grassland, Soil and Water Research Laboratory,
7 Temple, Texas, 2001.

8 Norbiato, D., Borga, M., Degli Esposti, S., Gaume, E., and Anquetin, S.: Flash flood warning
9 based on rainfall thresholds and soil moisture conditions: An assessment for gauged and
10 ungauged basins, *J. Hydrol.*, 362, 274-290, 2008.

11 Pan, M., and Wood, E.: Inverse streamflow routing, *Hydrol. Earth Syst. Sci. Discuss.*, 10,
12 6897-6929, 2013.

13 Parajka, J., Viglione, A., Rogger, M., Salinas, J. L., Sivapalan, M., and Blöschl, G.:
14 Comparative assessment of predictions in ungauged basins – Part 1:
15 Runoff-hydrograph studies, *Hydrol. Earth Syst. Sci.*, 17, 1783-1795,
16 doi:10.5194/hess-17-1783-2013, 2013.

17 Ponce, V., Hawkins, R., Golding, B., Smith, R., and Willeke, G.: Runoff curve number: Has it
18 reached maturity? *J. Hydrol. Eng.*, 1, 11-19, 1996.

19 Post, D. A., and Jakeman, A. J.: Predicting the daily streamflow of ungauged catchments in
20 SE Australia by regionalising the parameters of a lumped conceptual rainfall-runoff model,
21 *Ecol. Model.*, 123, 91-104, 1999.

22 Rallison, R. and Miller, N.: Past, present and future SCS runoff procedure, in: *Rainfall Runoff*
23 *Relationship*, edited by: Singh, V. P., Water Resour. Publ., Littleton, Colo., USA, 353–364,
24 1981.

25 Rakovec, O., Weerts, A., Hazenberg, P., Torfs, P., and Uijlenhoet, R.: State updating of a

1 distributed hydrological model with Ensemble Kalman Filtering: effects of updating
2 frequency and observation network density on forecast accuracy, *Hydrol. Earth Syst. Sci.*
3 *Discuss.*, 9, 3961–3999, 2012.

4 Reichle, R., McLaughlin, D., and Entekhabi, D.: Hydrologic data assimilation with the
5 ensemble Kalman filter, *Mon. Weather Rev.*, 130, 103-114, 2002.

6 Reichle, R. H., Crow, W. T., and Keppenne, C. L.: An adaptive ensemble Kalman filter for
7 soil moisture data assimilation, *Water Resour. Res.*, 44, W03423, 10.1029/2007wr006357,
8 2008.

9 Sellami, H., Jeunesse, I. L., Benabdallah, S., Baghdadi, N., and Vanclooster, M.: Uncertainty
10 analysis in model parameters regionalization: a case study involving the SWAT model in
11 Mediterranean catchments (Southern France), *Hydrol. Earth Syst. Sci. Discuss.*, 10,
12 4951-5011, 2013.

13 Sivapalan, M.: Prediction in ungauged basins: a grand challenge for theoretical hydrology,
14 *Hydrol. Process.*, 17, 3163-3170, 10.1002/hyp.5155, 2003.

15 Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Liang,
16 X., McDonnell, J. J., Mendiondo, E. M., O'Connell, P. E., Oki, T., Pomeroy, J. W., Schertzer,
17 D., Uhlenbrook, S., and Zehe, E.: IAHS Decade on Predictions in Ungauged Basins (PUB),
18 2003–2012: Shaping an exciting future for the hydrological sciences, *Hydrolog. Sci. J.*, 48,
19 857-880, doi:10.1623/hysj.48.6.857.51421, 2003.

20 Srinivasan, R., Zhang, X., and Arnold, J.: SWAT ungauged: hydrological budget and crop
21 yield predictions in the Upper Mississippi River Basin, *T. ASABE*, 53, 1533-1546, 2010.

22 Tran, A. P., Vanclooster, M., Zupanski, M., and Lambot, S.: Joint estimation of soil moisture
23 profile and hydraulic parameters by ground-penetrating radar data assimilation with
24 maximum likelihood ensemble filter, *Water Resour. Res.*, 50, 3131-3146,
25 10.1002/2013WR014583, 2014.

26 Troch, P. A., Paniconi, C., and McLaughlin, D.: Catchment-scale hydrological modeling and

1 data assimilation, *Adv. Water Resour.*, 26, 131-135, doi:10.1016/s0309-1708(02)00087-8,
2 2003.

3 Troch, P. A., Carrillo, G., Sivapalan, M., Wagener, T., and Sawicz, K.: Climate-vegetation-soil
4 interactions and long-term hydrologic partitioning: signatures of catchment co-evolution,
5 *Hydrol. Earth Syst. Sci.*, 10, 2927-2954, 2013.

6 van Griensven, A., Meixner, T., Grunwald, S., Bishop, T., Diluzio, M., and Srinivasan, R.: A
7 global sensitivity analysis tool for the parameters of multi-variable catchment models, *J.*
8 *Hydrol.*, 324, 10-23, doi:10.1016/j.jhydrol.2005.09.008, 2006.

9 Vrugt, J. A., Diks, C. G. H., Gupta, H. V., Bouten, W., and Verstraten, J. M.: Improved
10 treatment of uncertainty in hydrologic modeling: Combining the strengths of global
11 optimization and data assimilation, *Water Resour. Res.*, 41, W01017,
12 doi:10.1029/2004wr003059, 2005.

13 Vrugt, J. A., ter Braak, C. J., Diks, C. G., and Schoups, G.: Hydrologic data assimilation using
14 particle Markov chain Monte Carlo simulation: Theory, concepts and applications, *Adv. Water*
15 *Resour.*, 51, 457-478, 2013.

16 Vrugt, J. A., ter Braak, C. J. F., Clark, M. P., Hyman, J. M., and Robinson, B. A.: Treatment of
17 input uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain
18 Monte Carlo simulation, *Water Resour. Res.*, 44, W00b09, doi:10.1029/2007wr006720, 2008.

19 Wang, D., Chen, Y., and Cai, X.: State and parameter estimation of hydrologic models using
20 the constrained ensemble Kalman filter, *Water Resour. Res.*, 45, W11416,
21 doi:10.1029/2008wr007401, 2009.

22 Xie, X., and Zhang, D.: Data assimilation for distributed hydrological catchment modeling via
23 ensemble Kalman filter, *Adv. Water Resour.*, 33, 678-690,
24 doi:10.1016/j.advwatres.2010.03.012, 2010.

25 Xie, X., and Cui, Y.: Development and test of SWAT for modeling hydrological processes in
26 irrigation districts with paddy rice, *J. Hydrol.*, 396, 61-71, doi:10.1016/j.jhydrol.2010.10.032,

- 1 2011.
- 2 Xie, X.: Simultaneous State-Parameter Estimation for Hydrologic Modeling Using Ensemble
3 Kalman Filter. *Land Surface Observation, Modeling and Data Assimilation*, 441 - 464,
4 doi:10.1142/9789814472616_0014, 2013.
- 5 Xie, X., and Zhang, D.: A partitioned update scheme for state-parameter estimation of
6 distributed hydrologic models based on the ensemble Kalman filter, *Water Resour. Res.*, , 49,
7 7350 - 7365, doi:10.1002/2012WR012853, 2013.
- 8 Yang, J., Gong, P., Fu, R., Zhang, M., Chen, J., Liang, S., Xu, B., Shi, J., and Dickinson, R.:
9 The role of satellite remote sensing in climate change studies, *Nature Clim. Change*, 3, 875 -
10 883, doi:10.1038/nclimate1908, 2013.
- 11 Zhang, X., Srinivasan, R., and Van Liew, M.: Multi-site calibration of the SWAT model for
12 hydrologic modeling, *T. ASABE*, 51, 2039-2049, 2008.

1 Table 1. Model parameters to be estimated in data assimilation.

No.	Parameter (Type)	Description	Scale ⁽¹⁾	Process	Min	Max
1	CN ₂	SCS runoff curve number for moisture condition II (-)	HRU	Runoff	35.0	98.0
2	CH_K	Effective hydraulic conductivity of channels alluvium (mm/hour)	Subbasin	Channel water	0.02	76.0
3	SOL_AWC	Available water capacity of the soil layer (mm/mm soil)	HRU	Soil	0.0	1.0
4	SURLAG	Surface runoff lag coefficient (day)	HRU	Runoff	1.0	10.0
5	GWQMN	Threshold depth of water in the shallow aquifer required for return flow to occur (mm)	HRU	Groundwater	20.0	1000.0
6	ESCO	Plant evaporation compensation factor (-)	HRU	Evaporation	0.0	1.0
7	ALPHA_BF	Baseflow alpha factor (day)	HRU	Lateral water	0.0	1.0

2 * The hydrologic variables are with respect to the scales to reflect the related hydrologic
3 processes.

1 Table 2. Dynamic hydrologic states and outputs to be updated in data assimilation.

Variable	Description	Scale ⁽¹⁾	Storage ⁽²⁾
<i>Qsufstor</i>	Amount of surface runoff stored or lagged (mm)	HRU	QW
<i>Qlatstor</i>	Amount of lateral flow stored or lagged (mm)	HRU	QW
<i>Qpregw</i>	Amount of groundwater flow into the main channel (mm)	HRU	QW
<i>Wsol</i>	Amount of water stored in the soil layer for each HRU (mm)	HRU×Nlay	SW
<i>SM</i>	Amount of water stored in soil profile (mm)	Subbasin	SW
<i>Qshall</i>	Amount of shallow water stored or lagged (mm)	HRU	SW
<i>Qrchrg</i>	Amount of recharge entering the aquifer (mm)	HRU	SW
<i>Wr</i>	Amount of water stored in the reach (m ³)	Subbasin	CW
<i>Wb</i>	Amount of water stored in the bank (m ³)	Subbasin	CW
<i>Qr</i>	Amount of water flow out of reach (Streamflow, m ³ /s)	Subbasin	CW

2 ⁽¹⁾ This column indicates the scale at which each variable is simulated. Nlay is the number of
3 soil layers (Nlay = 4 for this study), and HRU×Nlay means the soil profile of each HRU is
4 partitioned into Nlay layers. ⁽²⁾ This column denotes water storage condition for each variable:
5 QW, quick water storage; SW, slow water storage; and CW, river channel storage.

6

1 Table 3. Fractional factors used to perturb the precipitation (η_p), simulated streamflow (η_{Q_m})
 2 and the observed streamflow (η_{Q_o}).

Distribution parameter	η_p	η_{Q_m}	η_{Q_o}
Values of fractional factor	0.10	0.15	0.10

3

1 **Figure captions**

2 Figure 1. Zhanghe River basin in China (a), the land use (b) and subbasin distribution with
3 DEM (C).

4 Figure 2. Streamflow prediction errors from the control-run simulation (left column) and the
5 data assimilation of scenario ASS_D (right column), i.e., only the observed streamflow from
6 Gauge D (outlet) is assimilated to update model states and parameters.

7 Figure 3. Streamflow prediction errors from scenarios ASS_BD and ASS_AB. Only the
8 results for Gauge C are shown because Gauge C is at the outlet of a pseudo-ungauged
9 subbasin in both scenarios.

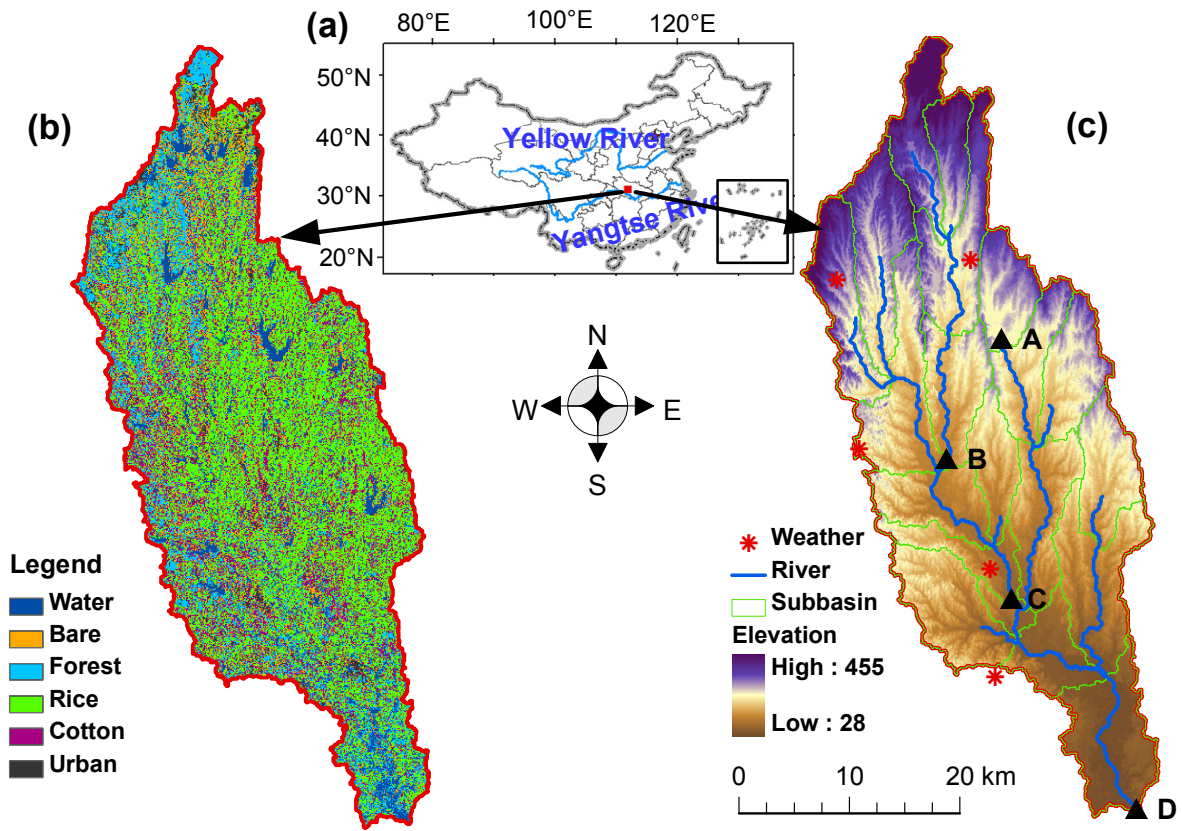
10 Figure 4. Estimations of two typical parameters (CN₂ and CH_K) from the ASS_D scenario.
11 The histograms in each plot, fitted with the Gaussian distribution function, represent the
12 ensemble distribution at three time steps.

13 Figure 5. Ensemble spreads (EnSp) of the seven parameters listed in Table 1:

14
$$EnSp = \sqrt{\frac{1}{Nu} \sum_{i=1}^{Nu} VAR_{En}(i)}$$
, where Nu is the number of HRUs or subbasins and $VAR_{En}(i)$

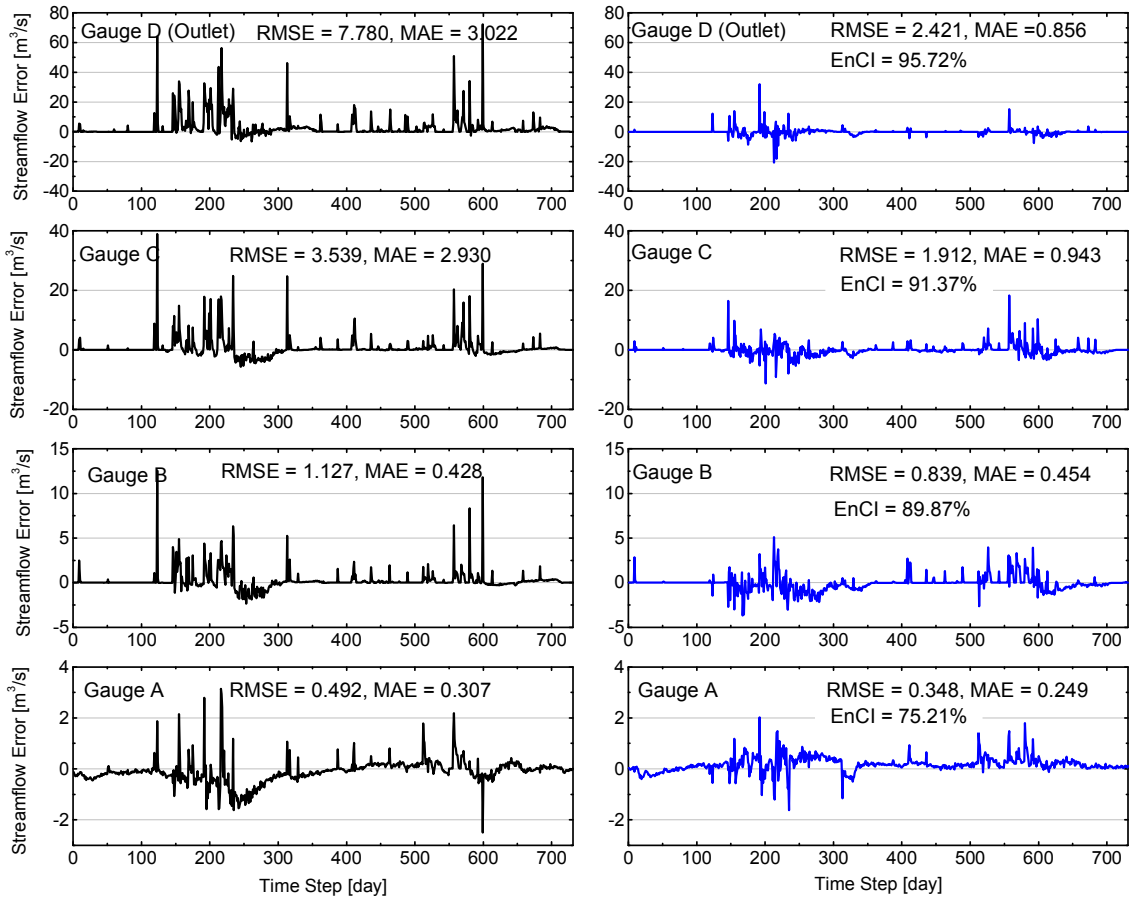
15 denotes the ensemble variance at each HRU or subbasin with respect to each parameter.

16 Figure 6. Streamflow predictions using four scenarios of different parameter sets. Only results
17 of Gauge C and D are shown.



1
2
3

Figure 1. Zhanghe River basin in China (a), the land use (b) and subbasin distribution with DEM (c).



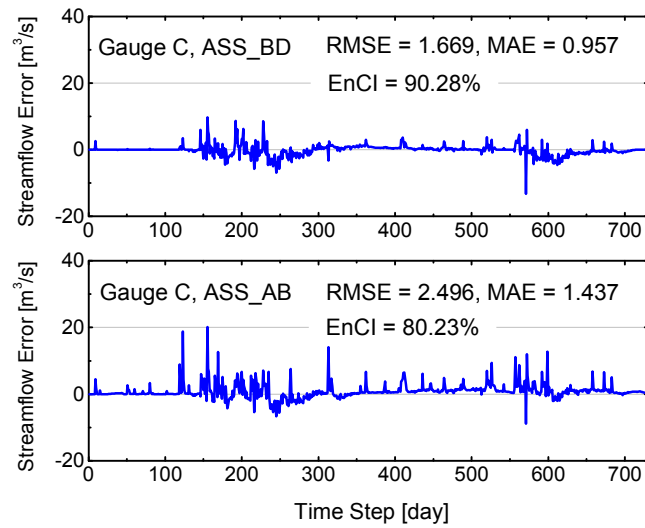
1

2

3

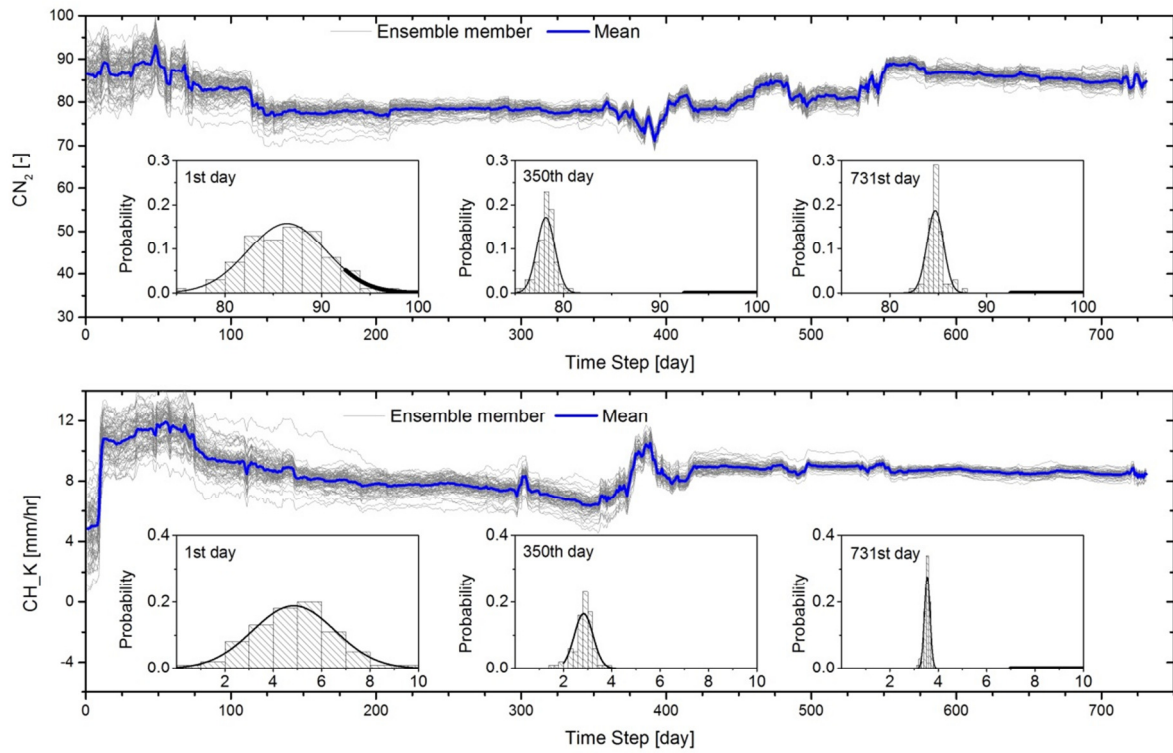
4

Figure 2. Streamflow prediction errors from the control-run simulation (left column) and the data assimilation of scenario ASS_D (right column), i.e., only the observed streamflow from Gauge D (outlet) is assimilated to update model states and parameters.



1

2 Figure 3. Streamflow prediction errors from scenarios ASS_BD and ASS_AB. Only the
 3 results for Gauge C are shown because Gauge C is at the outlet of a pseudo-ungauged
 4 subbasin in both scenarios.

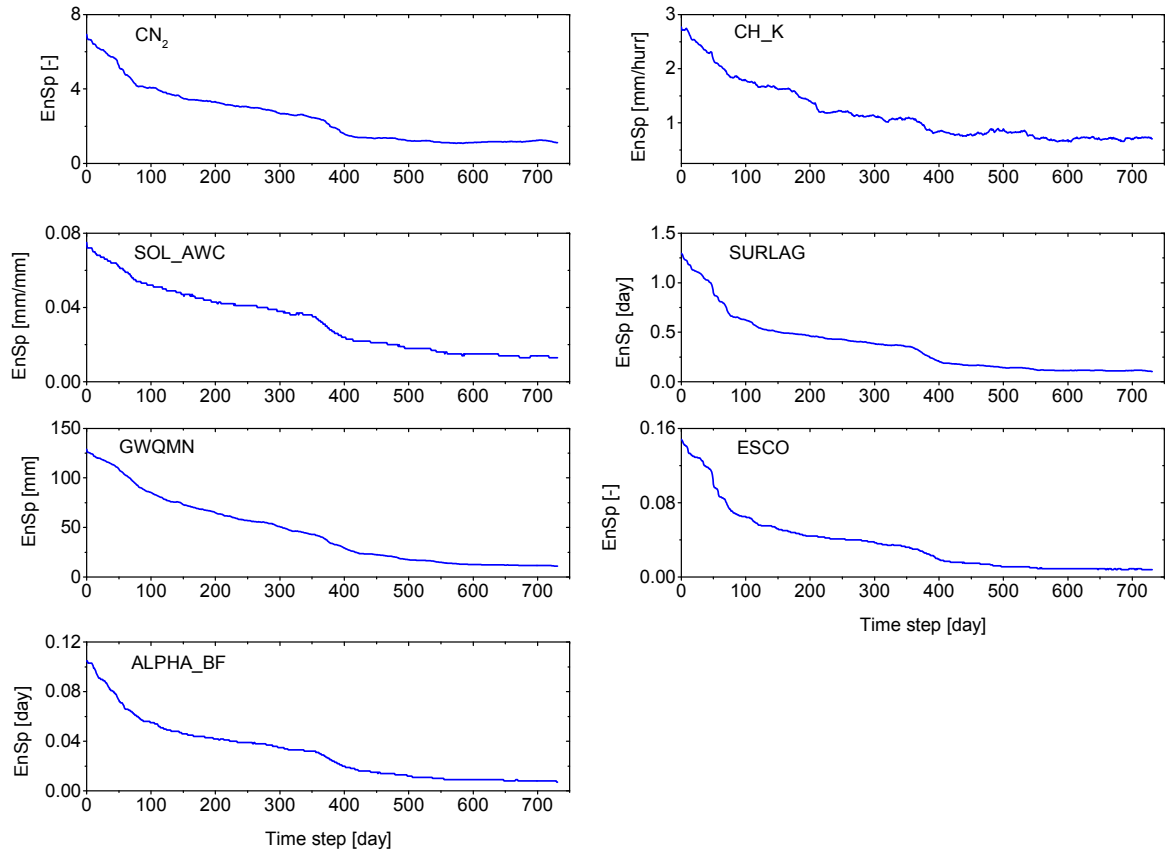


1

2 Figure 4. Estimations of two typical parameters (CN_2 and CH_K) from the ASS_D scenario.

3 The histograms in each plot, fitted with the Gaussian distribution function, represent the

4 ensemble distribution at three time steps.

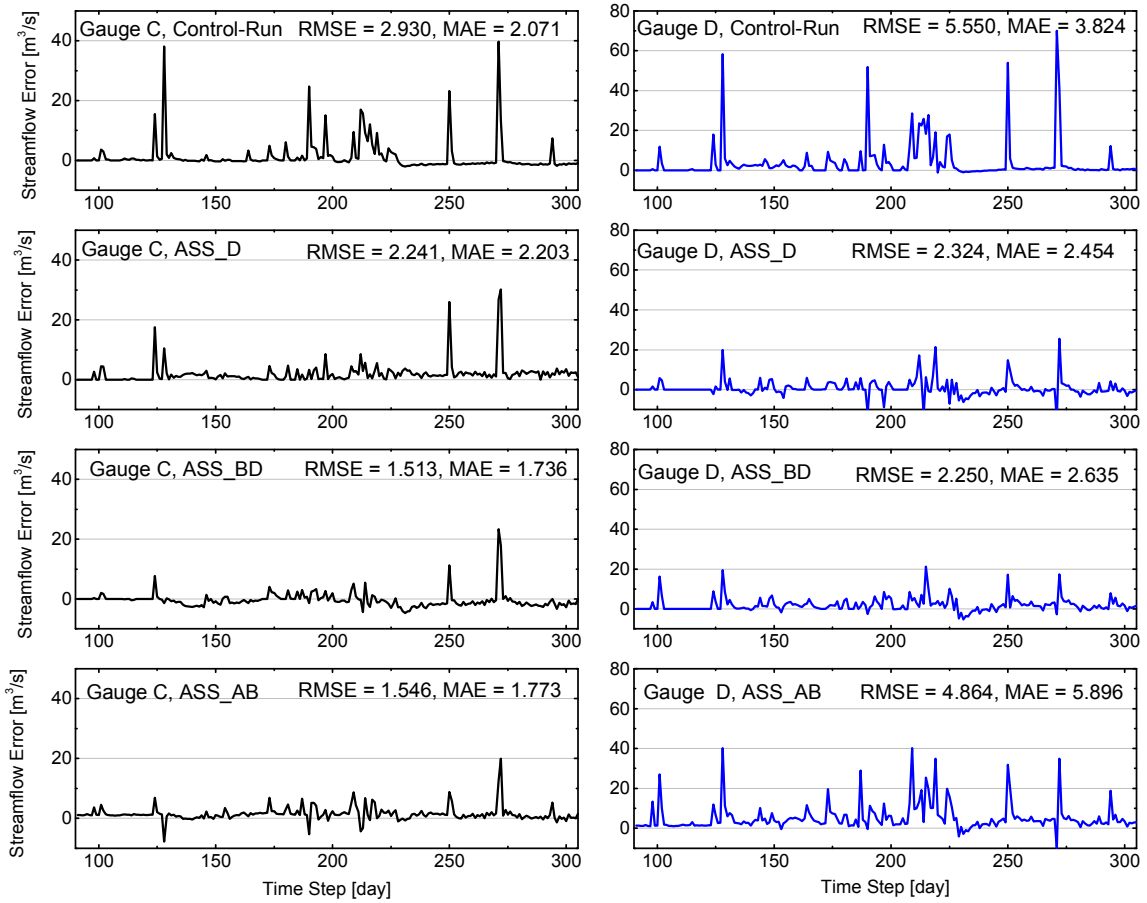


1

2 Figure 5. Ensemble spreads (EnSp) of the seven parameters listed in Table 1:

3
$$EnSp = \sqrt{\frac{1}{Nu} \sum_{i=1}^{Nu} VAR_{En}(i)}$$
, where Nu is the number of HRUs or subbasins and $VAR_{En}(i)$

4 denotes the ensemble variance at each HRU or subbasin with respect to each parameter.



1
 2 Figure 6. Streamflow predictions using four scenarios of different parameter sets. Only results
 3 of Gauge C and D are shown.