

Stochastic Environmental Research and Risk Assessment
Joint inversion of physical and geochemical parameters in groundwater models by sequential ensemble-based optimal design
--Manuscript Draft--

Manuscript Number:	SERR-D-17-00190R1	
Full Title:	Joint inversion of physical and geochemical parameters in groundwater models by sequential ensemble-based optimal design	
Article Type:	Original research	
Keywords:	Optimal sampling strategy; Physical and geochemical heterogeneity; Parameter estimation; Reactive transport model; Data assimilation.	
Corresponding Author:	Xiaoqing Shi Nanjing University Nanjing, Jiangsu CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Nanjing University	
Corresponding Author's Secondary Institution:		
First Author:	Tian Lan	
First Author Secondary Information:		
Order of Authors:	Tian Lan Xiaoqing Shi Beilei Jiang, Ph.D. Yuanyuan Sun, Ph.D. Jichun Wu, Ph.D.	
Order of Authors Secondary Information:		
Funding Information:	National Nature Science Foundation of China grant (U1503282)	Professor Jichun Wu
	National Nature Science Foundation of China grant (41672229)	Dr. Xiaoqing Shi
	National Nature Science Foundation of China grant (41172206)	Dr. Xiaoqing Shi
Abstract:	<p>Joint inversion of physical and geochemical parameters in groundwater reactive transport models is still a great challenge due to the intrinsic heterogeneities of natural porous media and the scarcity of observation data. In this study, we make use of a sequential ensemble-based optimal design (SEOD) method to jointly estimate physical and geochemical parameters of groundwater models. The effectiveness and efficiency of the SEOD method are illustrated by the comparison between the sequential optimization strategy and the conventional strategy (using fixed sampling locations) for two synthetic cases. Since the SEOD method is an optimization method based on the Ensemble Kalman Filter (EnKF), it invokes the time-consuming Genetic Algorithm (GA) at every assimilation step of the EnKF to obtain the optimal sampling locations. To enhance its computational efficiency, we improve the SEOD method by replacing the EnKF with the Ensemble Smoother with multiple data assimilation (ES-MDA). Furthermore, the influence factors of the original and improved SEOD method are also discussed. Our results show that the SEOD method provides an effective designed sampling strategy to accurately estimate heterogeneous distribution of physical and geochemical parameters. Moreover, the improved SEOD method is more</p>	

	advantageous than the original one in computational efficiency, making this SEOD framework more promising for future application.
Response to Reviewers:	More details please see the attached Response Letter.

Dear Editor,

We greatly appreciate your concern in our research paper. Your suggestions and both reviewers' comments contributed to its improvement. The reviewers' comments are all considered and the answers are given in the response letter. All the corrections are shown in the annotated copy of our manuscript.

Thank you for your consideration of this manuscript.

Sincerely,

Xiaoqing Shi

Reply to review comments on "Joint inversion of physical and geochemical parameters in groundwater models by sequential ensemble-based optimal design" (SERR-D-17-00190)

(Original review comments in black; **Reply from the authors in red**)

Comments from Chief Editor George Christakos:

Based on the advice received, I feel that your manuscript could be reconsidered for publication should you be prepared to incorporate major revisions. When preparing your revised manuscript, you are asked to carefully consider the reviewer comments which can be found below, and submit a list of responses to the comments. You are kindly requested to also check the website for reviewer attachment.

Response: We thank the editor for giving us the opportunity to address the AE's and reviewers' comments. We have carefully made the revisions corresponding to all the comments and provided a detailed point-by-point response as given below.

Comments from Associate Editor:

1- I received comments from two reviewers, reviewer 1 suggested revision and reviewer 2 suggested major revision. In general, reviewer 1 was positive to the manuscript, but asked the authors to provide more clear description of the SEOD and methodology used in this study. Reviewer 2, however, concerned about the novelty of this manuscript. The reviewer pointed out that the method used in the manuscript, SEOD method, was a known model, the authors need to highlight the innovation in the paper.

After receiving the comments, I carefully read the manuscript again and concur with the reviewers in general. I suggest major revision. The authors need to clearly highlight the manuscript's novelty and contribution. They need to reply the comments and question in detail and make revision to their manuscript accordingly.

Response: We would like to thank the AE for the constructive comments and for keeping the door for addressing reviewers' comments. Although the sequential ensemble-based optimal design (SEOD) method is a known method (developed by Man et al., 2016), the novelties and contributions of this study were highlighted in three aspects:

- (1) To our knowledge, the SEOD method is firstly applied for groundwater reactive transport modeling in this study. Major challenge for subsurface reactive transport models is the correct description of the subsurface processes using a limited amount of measurements. For this reason, to avoid complexity not justified by the available data, the spatial distributions of physical and geochemical parameters are often simplified. Man et al. (2016) developed the SEOD method, and demonstrated the effectiveness of this method by estimating only physical parameters in unsaturated flow models. In this study, we applied the SEOD method into jointly estimating physical and geochemical parameters of reactive transport model.
- (2) To distinguish this study from a mere case study, the worth of multiple observation data (piezometric head and concentration) was analyzed in the context of parameter estimation and prediction uncertainty reduction. The beneficial impact of assimilating dynamic multiple observation data in a transient flow and transport model is patent for the EnKF framework. Few studies have attempted to assimilate head and concentration to accurately estimate heterogeneous distribution of both physical and geochemical parameters. It is noted that the SEOD method was assimilated only head data in groundwater flow model as demonstrated in Man et al. (2016).
- (3) Additionally, we improved the SEOD method using the Ensemble Smoother with multiple data assimilation (ES-MDA) to replace the EnKF to enhance its computational efficiency. The comparison was illustrated in the section 4 of the revised manuscript.

According to the reviewer's constructive comments, we rewrote the Introduction, Abstract and Conclusion to highlight our contributions. To illustrate how we improved the SEOD framework using the ES-MDA, we also added details about the procedure of the SEOD method in section 2 of Methodologies.

2- In addition, the authors must link directly the paper's analysis to papers that have appeared in SERRA. This is because there is a question of relevancy, i.e., whether SERRA readers will benefit maximally by a paper that does not show strong evidence of relevancy to the journal's contents.

Response: To benefit maximally for SERRA readers, we added more literature reviews on widespread use of EnKF in parameter estimation and state prediction published in SERRA in section of "Introduction", for example: The Ensemble Kalman Filter (EnKF, Evensen 2003, 2009) is one of the most popular inverse methods over the last decade (Aanonsen et al. 2009; Oliver and Chen 2011; Zhou et al. 2014), recently used in parameter estimation and state prediction (Chen and Zhang 2006; Huang et al. 2009; Tong et al. 2010).

Comments from Reviewer #1:

This paper proposes to implement SEOD method into the groundwater transport modeling to improve the sampling strategy and parameter estimation process. The authors tested the SEOD method in two synthetic groundwater modeling cases and showed that it performs better than the traditional data assimilation method.

General comments:

1- I believe the manuscript is well written. I do like the objectives, and simple logic flow of this paper. The techniques can be helpful for other groundwater modelers.

Response: Thank you for the positive comment.

2- However, some points of this paper need to be clarified and more discussions are needed, as listed below. I would be in favor of publication after the authors addressed the comments given below.

Response: We have revised the manuscript carefully according to your comments and provided a detailed point-by-point response as given below.

Specific comments:

1. Since SEOD method is essential for this manuscript, I believe more details of this algorithm should be provided in the methodology section instead of just mentioning it (Page 8).

Response: Accepted and clarified. The SEOD method used in this study was adopted from previous work published in a recently WRR paper (Man et al., 2016). The details could be found in the corresponding cited literature and hence were only briefly discussed in the original manuscript. According to your comments, especially to help readers to understand how we improve the SEOD framework using the ES-MDA, the details about the SEOD method were added in section 2 Methodologies in the revised manuscript.

2. The methodology used in this research was not described clearly enough, please provide more details about the SEOD, EnKF, and relative entropy (RE).

Response: Accepted. More details about the SEOD method, the EnKF and the definition of RE were added in the revised manuscript (see section 2.2).

3. The novelty of this work is not well presented, since the methodology of this research is not new, please add more descriptions about what the new knowledge we can learn from this work (especially in the discussion and conclusion sections).

Response: Accepted and clarified. The Introduction, Discussion and Conclusion were modified to highlight the primary objective and key findings. This point has been addressed in our response to the AE's first comment (above). Although the SEOD method was a known method (developed by Man et al., 2016), the novelties and contributions of this study were highlighted in three aspects: (1) To our knowledge, the SEOD method is first applied for groundwater reactive transport modeling. In this study, we applied the SEOD method into jointly estimating physical and geochemical parameters of reactive transport model. (2) To distinguish this study from a mere case study, the worth of multiple observation data (piezometric head and concentration) was analyzed in the context of parameter estimation and prediction uncertainty reduction. The beneficial impact of assimilating dynamic multiple observation data in a transient flow and transport model is patent for the EnKF framework. Few studies have attempted to assimilate head and concentration to accurately estimate heterogeneous distribution of both physical and geochemical parameters. (3) Additionally, we improved the SEOD method using the ES-MDA to replace the EnKF to enhance its computational efficiency. The comparison was illustrated in the section 4 of the revised manuscript.

4. Please explain why use synthetic cases instead of realistic ones.

Response: Clarified. As we mentioned before, this study firstly explores an effective designed sampling strategy to accurately estimate heterogeneous distribution of physical and geochemical parameters in groundwater reactive transport models. Synthetic cases are good choices to illustrate the performance of a new proposed method due to the complexity and large computational burden of realistic ones. It is easily to extend our work to real practical applications in further study.

5. Page 6, line 8-9: please provide some new references about uncertainty study.

Response: Accepted. We have added several recently references, which relate to designing optimal sampling strategy in groundwater models for improving parameter identification and minimizing prediction uncertainty. Please see the section "Introduction" in the revised manuscript.

6. Please avoid using the terms like "firstly" in the manuscript, just "first".

Response: Accepted. We have checked and corrected in the revised manuscript.

7. The figure formats have some issues, please keep them consistent and put all the legends in the same place and add color bars for each line of subplots (i.e., Fig.8).

Response: Accepted. To our best, we have checked and corrected the figure formats, i.e., added color bars for each line, kept the consistent legends and put them in the same place. The legends were usually put in the upper right corner. However, for special cases that if there are important messages in the upper right corner of plots, we would change the location of legend correspondingly to make the plots more clearly for readers.

Comments from Reviewer #2:

1- The authors applied the ensemble Kalman filter (EnKF) method in joint inversion of physical and geochemical parameters of reactive transport model. And they used the sequential ensemble-based optimal design (SEOD) method to obtain the informative measurements at every assimilation step of the EnKF method. English is very easy to read. Authors have done much interesting work,

Response: Thank you for the positive comment.

2- however, SEOD method is a known model, authors do not describe its research situation to highlight their innovation in the paper.

Response: Accepted and clarified. This point has been addressed in our response to the AE and the first reviewer's comment 3 (above). The Introduction, Discussion and Conclusion sections were modified to highlight the innovation of this study in the revised manuscript.

3- Also they just present figures and analyze them simply, and the discussion section should be rewritten with more analysis. So I recommend major revision of this paper to the SERRA journal with its present format.

Response: Accepted. The content of discussion section was indeed not abundant enough in the original manuscript. In the revised manuscript, we focused on rewriting this section.

The original SEOD method is time-consuming because it invokes the Genetic Algorithm many times to estimate the optimal sampling design. To improve its computational efficiency, we replaced the EnKF with the ES-MDA. In discussion section, we compared the performance of the original and improved SEOD method. The effect of optimal sampling number and group division strategy on the result of this improved method were also discussed. The improvement increases the value of this method and the implications of this study.

1 **Joint inversion of physical and geochemical parameters in**
2 **groundwater models by sequential ensemble-based optimal design**

3 Tian Lan¹, Xiaoqing Shi^{1,*}, Beilei Jiang², Yuanyuan Sun¹, Jichun Wu^{1,*}

4

5 ¹ Key Laboratory of Surficial Geochemistry, Ministry of Education and School of
6 Earth Sciences and Engineering, Nanjing University, Nanjing, China 210023

7 ² Nanjing Hydraulic Research Institute, National Key Laboratory of Water Resources
8 and Hydraulic Engineering, Nanjing, China, 210029

9

10 *Corresponding author:

11 Xiaoqing Shi, Email: shixq@nju.edu.cn; (86) 25-89680839

12 Jichun Wu, Email: jcwu@nju.edu.cn; (86) 25-89680705

13

14 **Highlights:**

15 1. Both the heterogeneous physical and geochemical parameters of groundwater
16 models are accurately estimated by using the sequential ensemble-based optimal

17 design (SEOD) method.

18 2. The effectiveness and efficiency of the SEOD method are illustrated and

19 demonstrated by the comparison between the sequential optimization strategy

20 and the conventional strategy.

1 3. To enhance its computational efficiency, the SEOD method is improved by
2 replacing the EnKF with the Ensemble Smoother with multiple data assimilation
3 (ES-MDA).
4

1 **Abstract**

2 Joint inversion of physical and geochemical parameters in groundwater reactive
3 transport models is still a great challenge due to the intrinsic heterogeneities of natural
4 porous media and the scarcity of observation data. In this study, we make use of a
5 sequential ensemble-based optimal design (SEOD) method to jointly estimate
6 physical and geochemical parameters of groundwater models. The effectiveness and
7 efficiency of the SEOD method are illustrated by the comparison between the
8 sequential optimization strategy and the conventional strategy (using fixed sampling
9 locations) for two synthetic cases. Since the SEOD method is an optimization method
10 based on the Ensemble Kalman Filter (EnKF), it invokes the time-consuming Genetic
11 Algorithm (GA) at every assimilation step of the EnKF to obtain the optimal sampling
12 locations. To enhance its computational efficiency, we improve the SEOD method by
13 replacing the EnKF with the Ensemble Smoother with multiple data assimilation
14 (ES-MDA). Furthermore, the influence factors of the original and improved SEOD
15 method are also discussed. Our results show that the SEOD method provides an
16 effective designed sampling strategy to accurately estimate heterogeneous distribution
17 of physical and geochemical parameters. Moreover, the improved SEOD method is
18 more advantageous than the original one in computational efficiency, making this
19 SEOD framework more promising for future application.

20 **Keywords:** Optimal sampling strategy; Physical and geochemical heterogeneity;
21 Parameter estimation; Reactive transport model; Data assimilation.

1 **1 Introduction**

2 Joint inversion of physical and geochemical parameters in groundwater reactive
3 transport models is critical for reliable contaminant plume prediction, remediation and
4 management, but it is still a great challenge due to the intrinsic heterogeneities of
5 natural porous media and the scarcity of observation data. The subsurface
6 environment is highly variable in its physical and chemical composition.
7 Heterogeneity of physical parameters (e.g. hydraulic conductivity) has been shown to
8 exert a key control on the mixing and spreading of conservative solutes (Dagan 1984;
9 Rubin 1991; Sudicky 1986). For reactive solutes, their transport and reactions are
10 simultaneously influenced by geochemical parameters (Atchley et al. 2014; Li et al.
11 2010; Scheibe et al. 2006). Similar to physical heterogeneity, the heterogeneity of
12 geochemical parameters exists as well, which may be caused, for example, by spatial
13 variability in the activity of bacteria related to biodegradation (Fennell et al. 2001;
14 Sandrin et al. 2004). Therefore, it is important to jointly estimate the spatial
15 distribution of physical and geochemical parameters in groundwater reactive transport
16 models.

17 Inverse methods are often used by conditioning on observation data to
18 characterize the spatial variation of parameters, which has been extensively
19 investigated in the literature (e.g., Carrera et al. 2005; Dagan 1985; Doherty 2004;
20 Gómez-Hernández et al. 2003; Hendricks Franssen et al. 2009; Neuman 1973; Oliver
21 et al. 1997; Zhou et al. 2014). The Ensemble Kalman Filter (EnKF, Evensen 2003,
22 2009) is one of the most popular inverse methods over the last decade (Aanonsen et al.

1 2009; Oliver and Chen 2011; Zhou et al. 2014), recently used in parameter estimation
2 and state prediction (Chen and Zhang 2006; Huang et al. 2009; Tong et al. 2010). It is
3 a variant of the Kalman Filter (KF, Kalman 1960) based on the Monte Carlo method.
4 Unlike the KF, the EnKF was developed for nonlinear problems (Evensen 2003,
5 2009), its efficiency and effectiveness in nonlinear problems with high dimensionality
6 have been illustrated (Chen and Zhang 2006; Hendricks Franssen and Kinzelbach
7 2008; Moradkhani et al. 2005; Sorensen et al. 2004). In addition to the EnKF, the
8 Ensemble Smoother (ES, Van Leeuwen and Evensen 1996) and its iterative variants,
9 like the Ensemble Smoother with multiple data assimilation (ES-MDA, Emerick and
10 Reynolds 2013), are popular as well. Unlike the EnKF, the ES and the ES-MDA
11 perform global update rather than sequential update during the data assimilation,
12 avoiding restarting models again and again, so they are of more simplicity and
13 computational efficiency than the EnKF.

14 Much research has focused on developing better methods based on the EnKF to
15 broaden its implementation scale and improve its accuracy (Chen and Oliver 2010;
16 Emerick and Reynolds 2011; Gu and Oliver 2007; Li and Reynolds 2009), with the
17 sampling locations fixed during the data assimilation (called the conventional strategy
18 in the following discussion). However, it is intuitive that the data worth of
19 measurements is dramatically influenced by sampling locations, and the parameter
20 estimation result can be improved if the measurements are more informative even
21 though the number of sampling locations is the same. There has been much research
22 revealed the effect of sampling strategies on the parameter uncertainty and predictive

1 uncertainty in groundwater models (Carrera and Neuman 1986; Cleveland and Yeh
2 1990; Knopman and Voss 1987; Nowak et al. 2010; Sun and Yeh 2007; Ushijima and
3 Yeh 2015; Zhang et al. 2015). In view of these two aspects, Man et al. (2016)
4 integrated a sequential optimal design and the information theory into the EnKF
5 framework seamlessly to provide the most informative measurements for more
6 accurate parameter estimation, and proposed a sequential ensemble-based optimal
7 design (SEOD) method. Man et al. (2016) demonstrated the effectiveness of this
8 method by estimating only physical parameters in unsaturated flow models,
9 assimilating only piezometric head data. However, the SEOD method developed by
10 Man et al. (2016) invokes the optimization algorithm (the Genetic Algorithm) at each
11 assimilation step, so its computational efficiency is not very satisfying. Furthermore,
12 to the best of our knowledge, few studies have focused on joint inversion of physical
13 and geochemical parameters by assimilating multiple kinds of data.

14 The objective of this study is to estimate both physical and geochemical
15 parameters accurately in groundwater models by using the recent proposed SEOD
16 method, and to enhance the computational efficiency of the SEOD method by
17 replacing the EnKF with the ES-MDA. The rest of the paper is organized as follows.
18 In Section 2, the groundwater reactive transport model and the SEOD method are
19 described. In Section 3, synthetic one-dimensional and two-dimensional groundwater
20 reactive transport model cases are constructed to jointly estimate the physical and
21 geochemical parameters by using the SEOD method. In Section 4, the comparison
22 between the sequential optimization strategy and the conventional strategy, and the

1 effects of the ensemble size and the number of optimal sampling locations are
 2 discussed. Furthermore, we improve the SEOD method by replacing the EnKF with
 3 the ES-MDA to enhance its computational efficiency, and make a comparison of the
 4 original and the improved SEOD method in Section 4.4. Conclusions are summarized
 5 in Section 5.

6 **2 Methodologies**

7 **2.1 Groundwater reactive transport model**

8 In this work, transient flow is assumed, as the following governing equation
 9 (Bear 1972),

$$10 \quad \nabla \cdot (K \nabla H) + W = \mu_s \frac{\partial H}{\partial t} \quad (1)$$

11 where $\nabla \cdot$ is the divergence operator; ∇ is the gradient operator; K is the hydraulic
 12 conductivity [LT^{-1}]; H is the hydraulic head [L]; W is the volumetric injection
 13 (pumping) flow rate per unit volume of the aquifer [LT^{-1}]; μ_s is the specific storage of
 14 the aquifer [L^{-1}]; t is the time [T].

15 The governing equation for the transport and reactions of aqueous species is
 16 defined as (Zheng 2006; Prommer and Post 2010):

$$17 \quad \frac{\partial C_n}{\partial t} = \nabla \cdot (D \cdot \nabla C_n) - \nabla \cdot (v C_n) + r_{reac,n} + \frac{q_s}{\theta} C_n^s \quad (2)$$

18 where C_n is the aqueous concentration of the n th component [ML^{-3}]; t is the time [T];
 19 D is the diffusion coefficient [$L^2 T^{-1}$]; $v = (-K \nabla H) / \theta$ [$L^2 T^{-1}$]; $r_{reac,n}$ is the
 20 concentration change of the n th component caused by reactions; q_s is the volumetric
 21 flow rate per unit volume of the aquifer [T^{-1}]; θ is the effective porosity; and C_n^s is the
 22 concentration of the source or sink flux of the n th component [ML^{-3}].

1 Eq. (1) is solved by the numerical code MODFLOW-2000 (Harbaugh et al.
2 2000), and Eq. (2) is solved by the numerical code MT3DMS (Zheng 2006).

3 **2.2 Sequential ensemble-based optimal design (SEOD) method**

4 The sequential ensemble-based optimal design (SEOD) method is a new recently
5 proposed optimal method based on the EnKF (Man et al. 2016). At each recursive step,
6 the SEOD method provides an optimal sampling strategy, giving the maximum value
7 of information metric. Then, the analysis equation of the EnKF is used to update
8 estimated parameters by assimilating the most informative measurements, obtained
9 based on the optimal sampling strategy.

10 In this work, relative entropy (RE), also known as the Kullback-Leibler
11 divergence (Kullback 1977), is used to measure the information content of the
12 posterior probability density function (pdf) relative to the prior pdf. If these two
13 distributions are both n -dimensional Gaussian, RE between these two distributions is
14 defined as:

15
$$RE = J_b + [\ln \det(\mathbf{B}\mathbf{A}^{-1}) + \text{Tr}(\mathbf{A}\mathbf{B}^{-1}) - n] / 2 \quad (3)$$

16 where $J_b = (\mathbf{a} - \mathbf{b})^T \mathbf{B}^{-1} (\mathbf{a} - \mathbf{b}) / 2$ is the signal part of RE ; $\det(\cdot)$ denotes the
17 determinant, $\text{Tr}(\cdot)$ denotes the trace; \mathbf{a} and \mathbf{A} denote the mean and covariance matrix
18 of prior statistics respectively; \mathbf{b} and \mathbf{B} denote the mean and covariance matrix of
19 posterior statistics respectively.

20 The loop of the SEOD method for parameter estimation is briefly recalled. More
21 details can be found in Man et al. (2016). In the EnKF, all the parameters of interest \mathbf{p}
22 are augmented with state variables \mathbf{h} into a joint state vector $\mathbf{x} = [\mathbf{p} \ \mathbf{h}]^T$. Before the

1 forecast step, an ensemble of N_e realizations of parameters is generated.
 2 (I) Forecast step
 3 Rerun the forward model G from time 0 to time step $j+1$ with parameters updated
 4 at time step j (Eq. (4)).

$$5 \quad \mathbf{x}_{i,j+1}^f = G(\mathbf{x}_{i,j}^a), \quad i=1,2,\dots,N_e \quad (4)$$

6 In the above equation, i is the ensemble member index, j is the time step index,
 7 superscripts f and a denote forecast and analysis, respectively.

8 (II) Optimal design

9 Given a specific sampling strategy \mathbf{H} , the possible realizations of
 10 measurements can then be expressed as $\mathbf{d}_i = \mathbf{H}'\mathbf{x}_i^f + \xi_i$. With the realizations of
 11 measurements, the updated ensemble can be obtained from the Eq. (6). According to
 12 the prior and posterior statistics (mean and covariance), the information metrics RE of
 13 each candidate sampling strategies can be calculated. By comparing the RE values of
 14 different candidate sampling strategies, the optimal sampling design \mathbf{H}_{opt} can be
 15 determined by solving the following optimization problem (Eq. (5)) with the help of
 16 the Genetic Algorithm (GA, Whitley 1994).

$$17 \quad \mathbf{H}_{opt} = \arg \max RE(\mathbf{H}) \quad (5)$$

18 (III) Analysis step

19 After obtaining the optimal sampling strategy, the actual measurements \mathbf{d} can
 20 be obtained and used in the analysis step (Eq. (6)).

$$21 \quad \mathbf{x}_{i,j+1}^a = \mathbf{x}_{i,j+1}^f + \mathbf{C}_{YD}(\mathbf{C}_{DD} + \mathbf{C}_D)^{-1}(\mathbf{d}_{obs_{i,j+1}} - \mathbf{d}_{i,j+1}), \quad i=1,2,\dots,N_e \quad (6)$$

22 In the above equation, \mathbf{C}_{YD} is the cross-covariance matrix between the forecast

1 state and the predicted data, \mathbf{C}_{DD} is the covariance matrix of the predicted data, \mathbf{C}_D
 2 is the covariance matrix of the measurements error, \mathbf{d}_{obs} is the perturbed
 3 observations with noise of covariance \mathbf{C}_D , and \mathbf{d} is the predicted data.

4 After the analysis step, the updated ensemble X^a is obtained. Then, go back to
 5 step (I), the updated ensemble obtained this step is implemented for the next step.

6 To evaluate the performance of parameter estimation, two commonly used
 7 indicators, the *RMSE* and the *Ensemble Spread*, are defined as:

$$8 \quad RMSE = \sqrt{\frac{1}{N_m} \sum_{i=1}^{N_e} (\bar{Y}_i - Y_i)^2} \quad (7)$$

$$9 \quad Ensemble \ Spread (ES) = \sqrt{\frac{1}{N_m} \sum_{i=1}^{N_e} \text{var}(Y_i)} \quad (8)$$

10 where \bar{Y} and Y are the estimated and the reference field respectively; $\text{var}(Y)$ is
 11 the ensemble variance of the field; N_m is the total number of nodes in the study
 12 domain; N_e is the ensemble size; i is the node index. The *RMSE* measures the
 13 accuracy of the estimation, while the *Ensemble Spread* measures the uncertainty of
 14 the estimation.

15 **3 Case studies**

16 **3.1 Case 1: One-dimensional synthetic case**

17 In this case, a one-dimensional confined aquifer with a starting head of 100 m is
 18 constructed, in which saturated transient flow is assumed. As shown in Fig. 1(a), we
 19 choose the horizontal aquifer to be 5 m×150 m and the grid space to be 5 m both in
 20 horizontal x and y direction. Then, a Trichloroethylene (TCE) leaking area with an
 21 initial concentration of 1000 mg/L is introduced into the aquifer, and the degradation

1 of TCE is assumed to follow first-order kinetic reaction. Furthermore, an injection
 2 well and a pumping well are set upstream and downstream respectively, and all
 3 boundaries of the aquifer are assumed to be impermeable. In this case, the spatial
 4 distribution of the hydraulic conductivity (K) and the first-order rate constant (k_{TCE})
 5 (Fig. 1(b), (c)) are jointly estimated. At every assimilation step, 2 optimal sampling
 6 locations are selected from 30 candidate locations to provide the most informative
 7 measurements. More details are given in Table 1 and 2.

8 [Figure 1]

9 [Table 1]

10 [Table 2]

11 The log saturated hydraulic conductivity $Y_1 = \ln(K)$ and the first-order rate
 12 constant $Y_2 = k_{\text{TCE}}$ are assumed to be Gaussian distributed, with mean $\mu_{Y_1}=1$ and
 13 $\mu_{Y_2}=0.17$ and variance $\sigma_{Y_1}=1$ and $\sigma_{Y_2}=0.47$ respectively. Two arbitrary locations
 14 (x_1, y_1) and (x_2, y_2) in the random field are assumed to be correlated in the following
 15 form:

$$16 \quad C_Y(x_1, x_2) = C_Y(x_1, y_1; x_2, y_2) = \sigma^2 \exp\left[-\frac{|x_1 - x_2|}{\lambda_x} - \frac{|y_1 - y_2|}{\lambda_y}\right] \quad (9)$$

17 where the horizontal correlation length $\lambda_x= 5$ m, and the vertical correlation length λ_y
 18 = 30 m. Here we use the Karhunen-Loeve (K-L) expansion (Zhang and Lu 2004) to
 19 parameterize the random field so as to achieve the reference fields and initial
 20 ensemble members. The measurement errors of the head and concentration data are
 21 assumed to follow the standard normal distribution with the standard deviation of 0.01
 22 m and 10^{-6} mg/L respectively. Since the SEOD method is sequential, the uncertainty

1 changes as real-time measurements are assimilated, which leads to the optimal
2 sampling locations changing with time. The optimal sampling locations at 10
3 assimilation steps are shown in Fig. 2. It shows that the optimal sampling locations
4 change with the flow and concentration field so as to obtain the most informative
5 measurements. It is interesting to note that most optimal sampling locations are
6 located at the front of the contaminant plume, which means that these locations can
7 provide the most informative measurements.

8 [Figure 2]

9 In Fig. 3, we plot the curves of the ensemble mean and the standard deviation at
10 different assimilation steps. It shows that, for both Y_1 and Y_2 , the ensemble mean at
11 the final assimilation step is very close to the reference field. Furthermore, the
12 ensemble standard deviation is high at the early steps; however, it reduces
13 dramatically after assimilating the most informative measurements from the optimal
14 sampling locations. As shown in Fig. 4, the *RMSE* and the *Ensemble Spread* of Y_1 and
15 Y_2 decrease as the assimilation step increases, which also suggests that estimated
16 fields are close to their reference fields and of low uncertainty.

17 [Figure 3]

18 [Figure 4]

19 To further illustrate the accuracy and uncertainty of the estimation, we also
20 evaluate the performance of data match and model prediction. Considering the
21 limitation of space, two wells (marked with black circles in Fig. 1(a)) are selected
22 randomly to show the following evaluation results.

1 The initial and final ensembles of Y_1 and Y_2 are taken into the synthetic model
2 respectively to calculate the head and concentration data, which are then compared
3 with the real observations. After data assimilation, the data calculated by the final
4 ensemble become closer to the real observations, as shown in Fig. 5 and 6. Overall,
5 model uncertainty is significantly reduced after assimilating the most informative
6 measurements from the optimal sampling locations.

7 [Figure 5]

8 [Figure 6]

9 **3.2 Case 2: Two-dimensional synthetic case**

10 In this case, saturated transient flow is assumed in a two-dimensional confined
11 aquifer with a starting head of 50 m. As shown in Fig. 7, we choose the horizontal
12 aquifer to be 105 m×65 m and the grid space to be 5 m both in horizontal x and y
13 direction. Three TCE leaking sources with the constant injection flow of 80 m³/d and
14 the constant concentration of 50 mg/L per well are set upstream in the aquifer.
15 Furthermore, the liner sorption reaction of the TCE is considered in this case. Besides
16 three injection wells, two pumping wells with the constant pumping flow of 120 m³/d
17 per well are set downstream (Fig.7). In addition, all boundaries of the aquifer are
18 assumed to be impermeable. In this case, the hydraulic conductivity (K) and the liner
19 sorption constant (k_d) are assumed to be spatially heterogeneous (Fig.9). At every
20 assimilation step, 2 optimal sampling locations are selected from 77 candidate
21 locations (Fig.7) to provide the most informative measurements. More details are
22 given in Table 2 and 3.

1 [Figure 7]

2 [Table 3]

3 The log saturated hydraulic conductivity $Y_1 = \ln(K)$ is assumed to be Gaussian
4 distributed with mean $\mu_{Y_1}=1$ and variance $\sigma_{Y_1}=1$. The horizontal and vertical
5 correlation lengths of Y_1 are 40 m and 20 m respectively. With these statistics, the
6 reference field and initial ensemble members of Y_1 can be generated by the K-L
7 decomposition based on Eq. (9). For field Y_2 , it is assumed that there is a positive
8 correlation between $Y_2 = \ln(k_d)$ and Y_1 , i.e. $Y_2 = 0.5 \times Y_1 - 15.95$, on which the generation
9 of reference field Y_2 and its initial ensemble members are based. In addition, the
10 measurement errors of the head and concentration data are assumed to follow the
11 standard normal distribution with the standard deviation of 0.01 m and 10^{-6} mg/L
12 respectively.

13 In Fig. 8, we plot the optimal sampling locations at each assimilation step. It
14 shows the tendency that locations of large gradient are more likely to be selected as
15 the optimal sampling locations. Overall, it shows that the choice of the optimal
16 sampling locations at each assimilation step changes with the flow and concentration
17 field to obtain the most informative measurements.

18 [Figure 8]

19 The contour maps of the ensemble mean and the standard deviation at different
20 assimilation steps are plotted in Fig. 9. It shows that, for both Y_1 and Y_2 , the contour
21 maps of the ensemble mean exhibit a pattern very similar to the reference fields. Even
22 just after 4 assimilation steps, the contour maps of the ensemble mean recover the

1 major features of the reference fields of Y_1 and Y_2 . Furthermore, the ensemble
2 standard deviations reduce dramatically after assimilating the most informative
3 measurements from the optimal sampling locations, indicating that the optimal
4 sampling strategy does play a crucial role in the model inversion though only 2
5 sampling locations are selected at every assimilation step.

6 [Figure 9]

7 The *RMSE* and the *Ensemble Spread* of Y_1 and Y_2 are plotted in Fig. 10. It shows
8 that, these two indicators gradually decrease as assimilation step increases and finally
9 reach a low value, suggesting that the estimations of Y_1 and Y_2 in this case are accurate
10 and effective. Meanwhile, the difference between the *RMSE* and the *Ensemble Spread*
11 is small, indicating that the SEOD method estimates the uncertainty properly.

12 [Figure 10]

13 To evaluate the performance of data match and model prediction, the initial and
14 final ensembles of Y_1 and Y_2 are taken into the synthetic model respectively to
15 calculate the head and concentration data, which are then compared with the real
16 observations. It should be noted that only three wells (black circles in Fig. 7) are
17 selected randomly from the study domain to show the evaluation results due to the
18 space limitation. As shown in Fig. 11 and 12, the data calculated by the final ensemble
19 are very close to the real observations, performing much better than those calculated
20 by the initial ensemble. It indicates that the estimated fields of Y_1 and Y_2 are both of
21 low uncertainty after assimilating the most informative measurements from the
22 optimal sampling locations.

1 [Figure 11]

2 [Figure 12]

3 **4. Discussion**

4 **4.1 Comparison of sequential optimization strategy and conventional strategy**

5 In order to illustrate and demonstrate the effectiveness and efficiency of the
6 SEOD method for jointly estimating physical and geochemical parameters, the
7 sequential optimization strategy is compared with the conventional strategy in this
8 subsection. For the convenience of comparison, several synthetic cases (Case 11, 12
9 13) are constructed based on Case 2 by replacing the sequential optimization strategy
10 with the conventional strategy (different fixed sampling locations numbers for
11 different cases). Except this, the other model parameters of cases constructed here are
12 the same as those of Case 2. More details are given in Table 2.

13 Fig.13 shows the *RMSE* and the *Ensemble Spread* for different cases. It
14 illustrates that the sequential optimization strategy obtains better performance of the
15 parameter estimation when the number of sampling locations is the same. Even more,
16 the sequential optimization strategy with 2 optimal sampling locations (Case 2)
17 performs better than the conventional strategy with 10 fixed sampling locations (Case
18 12). Besides, the conventional strategy with a large number of fixed sampling
19 locations could result in the *Ensemble Spread* becoming very small at the first few
20 assimilation steps, which could prevent assimilating further measurements.

21 [Figure 13]

1 **4.2 Effect of ensemble size**

2 All results shown so far of Case 2 are based on an ensemble of 100 realizations.

3 To evaluate the impact of the ensemble size on the parameter estimation, an analysis

4 with an ensemble of 50, 300, 500, 1000 realizations (Table 2) is performed here.

5 The *RMSE* and the *Ensemble Spread* of different cases are shown in Fig. 14

6 below. It shows that an appropriate ensemble size is important for the parameter

7 estimation. If the ensemble size is too small (Case 3), ensemble collapse, a

8 phenomenon in which the *Ensemble Spread* is artificially small relative to its *RMSE*,

9 could happen. If the ensemble size is too large (Case 5, 6), it could lead to more

10 computational burden and introduce more observation errors into the model as the

11 SEOD method is based on the Monte Carlo method. It shows that the *RMSE* and the

12 *Ensemble Spread* of Y_1 and Y_2 are small and close to each other when the ensemble

13 size is 100, suggesting that the estimations of Y_1 and Y_2 are accurate and the model

14 uncertainty is estimated properly. Accordingly, the ensemble size is set to 100 in the

15 cases discussed below.

16 [Figure 14]

17 **4.3 Effect of the number of optimal sampling locations**

18 Optimizing too many sampling locations could bring a heavy computational

19 burden. Here, to explore the impact of the number of optimal sampling locations on

20 the parameter estimation, several synthetic cases with different numbers of optimal

21 sampling locations are constructed. More details are given in Table 2.

22 As shown in Fig. 15 below, the *RMSE* is no longer sensitive to the number of

1 optimal sampling locations when the number of optimal sampling locations is large
2 enough, suggesting that there could be a threshold value of the number of optimal
3 sampling locations in this synthetic model. On the one hand, if the number of optimal
4 sampling locations is too large, the *Ensemble Spread* becomes extremely small at the
5 first few assimilation steps, which could prevent assimilating further measurements
6 into the model. On the other hand, too many optimal sampling locations could lead to
7 high economic cost and heavy computational burden. Therefore, 2 to 5 optimal
8 sampling locations are enough and appropriate in this model.

9 [Figure 15]

10 **4.4 Improvement of the SEOD method**

11 In the original SEOD method, since the EnKF is a sequential history matching
12 method, the optimization algorithm part (GA) needs to be invoked N_s times to obtain
13 the optimal sampling design at each assimilation step, which is time-consuming. To
14 enhance its computational efficiency, we improve the original SEOD method by
15 replacing the EnKF with the ES-MDA. The loop of the improved SEOD method is
16 shown in Fig. 16.

17 [Figure 16]

18 Fig. 16 shows that the loop of the improved SEOD method is divided into outer
19 loop and inner loop. The outer loop is similar to the original SEOD method, which
20 consists of a forecast step, an optimal design step and an analysis step. The inner loop
21 of the improved SEOD method is part of the ES-MDA. Unlike the EnKF, the
22 ES-MDA performs N_a times global update so as to assimilate the same data (all

1 available data) multiple times without restarting the forward model, which helps
 2 enhance the computational efficiency. In the improved SEOD method, we divide all
 3 N_s assimilation steps in the original SEOD method into N_g groups, the following loop
 4 is performed for each group in chronological order (from 1 to N_g). Note that N_a of all
 5 cases in this subsection (Case 14, 15, 16, 17) is set to 4 (Emerick and Reynolds 2013).

6 More details are given in Table 2.

7 (I) Forecast step

8 Run the forward model G from beginning time step of the group $j+1$ to the end
 9 time step of the group $j+1$ with updated parameters from the group j (Eq. (10)).

$$10 \quad \mathbf{x}_{i,j+1}^f = G(\mathbf{x}_{i,j}^a), \quad i=1,2,\dots,N_e \quad (10)$$

11 In the above equation, i is the ensemble member index, j is the group index (from
 12 1 to N_g), superscripts f and a denote forecast and analysis, respectively.

13 (II) Optimal design

14 This step is similar with the original SEOD method, using the GA to solve an
 15 optimization problem to obtain the most informative measurements from the optimal
 16 sampling design.

17 (III) Analysis step

18 The following update equation (Eq. (11)) of the ES-MDA is different from that
 19 of the EnKF.

$$20 \quad \mathbf{x}_{i,j+1}^a = \mathbf{x}_{i,j+1}^f + \mathbf{C}_{YD} (\mathbf{C}_{DD} + \alpha_l \mathbf{C}_D)^{-1} (\mathbf{d}_{obs_{i,j+1}} - \mathbf{d}_{i,j+1}), \quad i=1,2,\dots,N_e \quad (11)$$

21 In the above equation, l is the times index of the ES-MDA, $l=1,2,\dots,N_a$; \mathbf{d}_{obs}
 22 is the perturbed observations with noise of covariance $\alpha_l \mathbf{C}_D$ ($\alpha_1=9.333$, $\alpha_2=7.0$,

1 $\alpha_3=4.0$ and $\alpha_4=2.0$, Emerick and Reynolds 2013). Other letters in this equation have
2 the same meaning as those in Eq. (6).

3 After N_a times global update, the updated ensemble X^a of the group $j+1$ is
4 obtained here. Then, go back to step (I), the updated ensemble is implemented for the
5 next group. Through this improvement, the times of invoking the GA decrease from
6 N_s to N_g , which helps enhance the computational efficiency.

7 To compare the improved SEOD method with the original one and discuss the
8 influence factors of the improved one, several cases are constructed with all model
9 parameters the same as those in Case 2. More details can be found in Table 2. The
10 results of these cases are shown in Fig. 17.

11 [Figure 17]

12 Fig. 17 shows that the number of optimal sampling locations dominantly affects
13 the results of the improved SEOD method. In Case 2, two optimal sampling locations
14 are chosen at each assimilation step. During the whole data assimilation, total 20
15 sampling locations are used to obtain measurements at most (2×10 , if the optimal
16 sampling locations are different at each step). Therefore, if the number of optimal
17 sampling locations in the improved SEOD method is too small, the improved SEOD
18 method can't estimate the parameters of the entire study domain well just through a
19 few sampling locations (e.g., Case 14 with total 4 optimal sampling locations at most).

20 When the number of sampling locations is enough, the result of the improved SEOD
21 method is acceptable (e.g., Case 15 with total 10 optimal sampling locations at most).
22 The result of Case 15 is comparable with the result of Case 2, but the computer cost of

1 Case 15 is much less than that of Case 2 because that Case 15 just invokes the GA
2 only twice while Case 2 invokes the GA 10 times. Therefore, if the number of optimal
3 sampling locations is not set to a very small value, the computational efficiency will
4 be enhanced by using the improved SEOD method.

5 Furthermore, the way of dividing the observation time (assimilation steps in the
6 original SEOD method) into several groups (called the group division strategy in the
7 following context) affects the results of the improved method as well. From the
8 comparison of Case 15, Case 16 and Case 17, it is obvious that Case 16, whose
9 observation time in each divided groups is progressively increasing, has a better data
10 assimilation result than the other two cases, whose observation time in each divided
11 groups is equivalent or progressively decreasing. It is an interesting phenomenon,
12 which is worth further research. For now, we think it is probably because more and
13 more precise measurements in the early stage of the data assimilation would lead to
14 excessive update of estimated parameters (Burgers et al. 1998; Evensen 2009).
15 Therefore, future studies should focus on optimizing the group division strategy to
16 obtain a more accurate estimation of model parameters.

17 **5. Conclusions**

18 In this study, we make use of a sequential ensemble-based optimal design (SEOD)
19 method to jointly estimate physical and geochemical parameters of groundwater
20 models.

21 Both physical and geochemical parameters are estimated accurately in the
22 one-dimensional and two-dimensional synthetic cases by using the SEOD method.

1 Uncertainties of both physical and geochemical parameters decrease after assimilating
2 the most informative measurements at the optimal sampling locations, and the
3 accuracy of model prediction increase meanwhile. Furthermore, several comparison
4 cases are tested and analyzed, results illustrate and demonstrate the effectiveness and
5 efficiency of the SEOD method on jointly estimating high-dimensional physical and
6 geochemical parameters in groundwater models.

7 The ensemble size and the number of optimal sampling locations have impacts
8 on the parameter estimation based on the SEOD method. A too small ensemble size
9 would lead to the ensemble collapse. Furthermore, when the number of optimal
10 sampling locations is too large, heavier computational burden and more observation
11 errors would be caused, and the *RMSE* is no longer sensitive to the number of optimal
12 sampling locations. How to determine the optimal ensemble size and sampling
13 locations number for different scenarios is worth further investigation.

14 The original SEOD method has a heavy computational burden because it invokes
15 the GA too many times. To enhance its computational efficiency, we proposed an
16 improved SEOD method in this study by replacing the EnKF with the ES-MDA. The
17 results of comparison cases show that the improved SEOD method is advantageous
18 than the original one, which makes the SEOD framework more promising for the
19 parameter estimation and the optimal sampling strategy design. The number of
20 optimal sampling locations and the strategy of dividing groups would affect the
21 results of the improved SEOD method.

22 It is noted that only two kinds of measurements (head and concentration) are

1 assimilated in this work. More kinds of measurements (e.g., hydraulic conductivity,
2 porosity, temperatures and hydrogeophysical data) can be assimilated simultaneously
3 so as to make use of more hard and soft data to improve the accuracy of parameter
4 estimation in further study.

5

6 **Acknowledgements**

7 The authors would like to thank the anonymous referees for their insightful comments and
8 suggestions that have helped improve the paper. This work was financially supported by
9 the National Nature Science Foundation of China grants (No. U1503282, 41672229,
10 and 41172206). We would like to thank Mr. Jun Man from Zhejiang University for
11 providing the SEOD code.

12

13 **References**

- 14 Aanonsen SI, Nævdal G, Oliver DS, Reynolds AC, Vallès B (2009) The ensemble Kalman filter in
15 reservoir engineering--a review. *Spe Journal* 14:393-412
- 16 Atchley AL, Navarre-Sitchler AK, Maxwell RM (2014) The effects of physical and geochemical
17 heterogeneities on hydro-geochemical transport and effective reaction rates. *J Contam
Hydrol* 165:53-64 doi:10.1016/j.jconhyd.2014.07.008
- 18 Bear J (1972) Dynamics of Fluids in Porous Materials. Dover, New York
- 19 Burgers G, Leeuwen P, Evensen G (1998) Analysis scheme in the ensemble Kalman filter. *Month
Weather Rev* 126:1719–1724
- 20 Carrera J, Alcolea A, Medina A, Hidalgo J, Slooten LJ (2005) Inverse problem in hydrogeology.
21 *Hydrogeology Journal* 13(1):206–222
- 22 Carrera J, Neuman SP (1986) Estimation of Aquifer Parameters Under Transient and Steady State
23 Conditions: 3. Application to Synthetic and Field Data. *Water Resources Research*
24 22:228-242
- 25 Chen Y, Oliver DS (2010) Cross-covariances and localization for EnKF in multiphase flow data
assimilation. *Computational Geosciences* 14:579-601
- 26 Chen Y, Zhang D (2006) Data assimilation for transient flow in geologic formations via ensemble
27 Kalman filter. *Advances in Water Resources* 29:1107-1122
28 doi:10.1016/j.advwatres.2005.09.007

- 1 Cleveland TG, Yeh WWG (1990) Sampling network design for transport parameter identification.
2 Journal of Water Resources Planning & Management 116:764-783
- 3 Dagan G (1984) Solute transport in heterogeneous porous formations. Journal of Fluid Mechanics
4 145:151 doi:10.1017/s0022112084002858
- 5 Dagan G (1985) Stochastic modeling of groundwater flow by unconditional and conditional
6 probabilities: The inverse problem. Water Resources Research 21(1):65–72
- 7 Doherty J (2004) PEST: Model-Independent Parameter Estimation,User's Manual (5th edition).
8 Watermark Numerical Computing, Australia
- 9 Emerick AA, Reynolds AC (2011) Combining sensitivities and prior information for covariance
10 localization in the ensemble Kalman filter for petroleum reservoir applications.
11 Computational Geosciences 15:251-269
- 12 Emerick AA, Reynolds AC (2013) Ensemble smoother with multiple data assimilation. Comput
13 Geosci-Uk 55:3-15 doi:10.1016/j.cageo.2012.03.011
- 14 Evensen G (2003) The Ensemble Kalman Filter: theoretical formulation and practical
15 implementation. Ocean Dynamics 53:343-367
- 16 Evensen G (2009) Data assimilation: the ensemble Kalman filter. Springer, Berlin
- 17 Fennell DE, Carroll AB, Gossett JM, Zinder SH (2001) Assessment of indigenous reductive
18 dechlorinating potential at a TCE-contaminated site using microcosms, polymerase chain
19 reaction analysis, and site data. Environmental Science & Technology 35:1830-1839
- 20 Gómez-Hernández JJ, Hendricks Franssen HJ, Sahuquillo A (2003) Stochastic conditional inverse
21 modeling of subsurface mass transport: A brief review and the self-calibrating method.
22 Stochastic Environmental Research and Risk Assessment 17(5):319–328
- 23 Gu Y, Oliver DS (2007) An Iterative Ensemble Kalman Filter for Multiphase Fluid Flow Data
24 Assimilation. Spe Journal 12:1990 - 1995
- 25 Harbaugh AW, Banta ER, Hill MC, McDonald MG (2000) MODFLOW-2000, The U. S.
26 Geological Survey Modular Ground-Water Model—User Guide to Modularization
27 Concepts and the Ground-Water Flow Process. U.S. Geological Survey Open-File Report
28 00-92, 121 p
- 29 Hendricks Franssen HJ, Alcolea A, Riva M, Bakr M, van der Wiel N, Stauffer F, Guadagnini A
30 (2009) A comparison of seven methods for the inverse modelling of groundwater flow.
31 Application to the characterisation of well catchments. Advances in Water Resources
32 32(6):851–872
- 33 Hendricks Franssen HJ, Kinzelbach W (2008) Real-time groundwater flow modeling with the
34 Ensemble Kalman Filter: Joint estimation of states and parameters and the filter
35 inbreeding problem. Water Resour Res 44:354-358
- 36 Huang C, Hu BX, Li X, Ye M (2009) Using data assimilation method to calibrate a heterogeneous
37 conductivity field and improve solute transport prediction with an unknown
38 contamination source. Stochastic Environmental Research and Risk Assessment
39 23(8):1155
- 40 Kalman RE (1960) A new approach to linear filtering and prediction problems. Trans ASME J
41 Basic Eng 82(D):35–45
- 42 Knopman DS, Voss CI (1987) Behavior of sensitivities in the one-dimensional
43 advection-dispersion equation: Implications for parameter estimation and sampling design.
44 Water Resources Research 23:253-272

- 1 Kullback S (1997) Information theory and statistics. Courier Corporation
- 2 Li G, Reynolds AC (2009) Iterative Ensemble Kalman Filters for Data Assimilation. *Spe Journal*
14:496-505
- 3 Li L, Steefel CI, Kowalsky MB, Englert A, Hubbard SS (2010) Effects of physical and
geochemical heterogeneities on mineral transformation and biomass accumulation during
biostimulation experiments at Rifle, Colorado. *J Contam Hydrol* 112:45-63
doi:10.1016/j.jconhyd.2009.10.006
- 4 Man J, Zhang J, Li W, Zeng L, Wu L (2016) Sequential ensemble-based optimal design for
parameter estimation. *Water Resour Res* 52:7577-7592 doi:10.1002/2016wr018736
- 5 Moradkhani H, Sorooshian S, Gupta HV, Houser PR (2005) Dual state-parameter estimation of
hydrological models using ensemble Kalman filter. *Advances in Water Resources*
28:135-147
- 6 Neuman SP (1973) Calibration of distributed parameter groundwater flow models viewed as a
multiple objective decision process under uncertainty. *Water Resources Research* 9(4):1006–1021
- 7 Nowak W, De Barros FPJ, Rubin Y (2010) Bayesian geostatistical design: Task-driven optimal site
investigation when the geostatistical model is uncertain. *Water Resources Research* 46(3):
374-381 doi:10.1029/2009WR008312
- 8 Oliver DS, Chen Y (2011) Recent progress on reservoir history matching: a review. *Computational
Geosciences* 15(1):185-221
- 9 Oliver DS, Cunha LB, Reynolds AC (1997) Markov chain Monte Carlo methods for conditioning
a permeability field to pressure data. *Mathematical Geology* 29(1): 61–91
- 10 Prommer CH, Post V (2010) A Reactive Multicomponent Transport Model for Saturated Porous
Media. *Groundwater* 48(5):627-632
- 11 Rubin Y (1991) Transport in heterogeneous porous media: Prediction and uncertainty. *Water
Resour Res* 27:1723-1738
- 12 Sandrin SK, Brusseau ML, Piatt JJ, Bodour AA, Blanford WJ, Nelson NT (2004) Spatial
variability of in situ microbial activity: biotracer tests. *Groundwater* 42:374-383
- 13 Scheibe TD, Fang Y, Murray CJ, Roden EE, Chen J, Chien YJ, Brooks SC, Hubbard SS (2006)
Transport and biogeochemical reaction of metals in a physically and chemically
heterogeneous aquifer. *Geosphere* 2(4):220-235 doi:10.1130/Ges00029.1
- 14 Sorensen JVT, Madsen H, Madsen H (2004) Data assimilation in hydrodynamic modelling: on the
treatment of non-linearity and bias. *Stoch Environ Res Risk Assess* 18(7):228–244
- 15 Sudicky EA (1986) A natural gradient experiment on solute transport in a sand aquifer: Spatial
variability of hydraulic conductivity and its role in the dispersion process. *Water Resour
Res* 22:2069-2082 doi:10.1029/WR022i013p02069
- 16 Sun NZ, Yeh WWG (2007) Development of objective-oriented groundwater models: 2. Robust
experimental design. *Water resources research* 43(2) doi:10.1029/2006WR004888
- 17 Tong J, Hu BX, Yang J (2010) Using data assimilation method to calibrate a heterogeneous
conductivity field conditioning on transient flow test data. *Stochastic environmental
research and risk assessment* 24(8):1211-23
- 18 Ushijima TT, Yeh WWG (2015) Experimental design for estimating unknown hydraulic
conductivity in an aquifer using a genetic algorithm and reduced order model. *Advances
in Water Resources* 86:193-208

- 1 Van Leeuwen PJ, Evensen G (1996) Data Assimilation and Inverse Methods in Terms of a
2 Probabilistic Formulation. *Monthly Weather Review* 124:2898-2913
- 3 Whitley D (1994) A genetic algorithm tutorial. *Stat Comput* 4(2):65-85
- 4 Zhang D, Lu Z (2004) An efficient, high-order perturbation approach for flow in random porous
5 media via Karhunen-Loëve and polynomial expansions. *Journal of Computational
6 Physics* 194:773-794
- 7 Zhang J, Zeng L, Chen C, Chen D, Wu L (2015) Efficient Bayesian experimental design for
8 contaminant source identification. *Water Resources Research* 51(1):576-598
- 9 Zheng C (2006) MT3DMS v5.2 supplemental user's guide: Technical report to the U.S. Army
10 Engineer Research and Development Center, Department of Geological Sciences,
11 University of Alabama, p 24
- 12 Zhou HY, Gómez-Hernández JJ, Li LP (2014) Inverse methods in hydrogeology: Evolution and
13 recent trends. *Advances in Water Resources* 63:22-37
- 14

Tables

Table 1: Flow and transport parameters used in Case 1

Flow simulation	Transient state
Total simulation time (days)	10
Stress period	1
Time steps	100
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	150
Model width (m)	5
Model height (m)	5
Starting head (m)	100
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
Injection rate per well (m^3/d)	50
Pumping rate per well (m^3/d)	45

Table 2: Data assimilation related parameters used in different cases

Case name	Dimension	Number of ensemble (N_e)	Number of assimilation step (N_s)	Optimize or not	Number of optimal sampling locations
Case 1	1	300	10	Y	2
Case 2	2	100	10	Y	2
Case 3	2	50	10	Y	2
Case 4	2	300	10	Y	2
Case 5	2	500	10	Y	2
Case 6	2	1000	10	Y	2
Case 7	2	100	10	Y	1
Case 8	2	100	10	Y	5
Case 9	2	100	10	Y	10
Case 10	2	100	10	Y	20
Case 11	2	100	10	N	(2 fixed)
Case 12	2	100	10	N	(10 fixed)
Case 13	2	100	10	N	(20 fixed)
Case 14	2	100	8(50, 50)*	Y	2
Case 15	2	100	8(50, 50)	Y	5
Case 16	2	100	12(20, 30, 50)	Y	5
Case 17	2	100	12(50, 30, 20)	Y	5

2 * The numbers in the parentheses are the group division of observation time, and the number in front of
3 the parentheses is the number of assimilation steps. For example, 8(50, 50) represents that there are 8
4 steps in the assimilation and the observation time is divided into two groups with each group having an
5 observation time of 50 days.

1 **Table 3: Flow and transport parameters used in Case 2**

Flow simulation	Transient state
Total simulation time (days)	100
Stress period	1
Time steps	200
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	105
Model width (m)	65
Model height (m)	5
Starting head (m)	50
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
horizontal transverse dispersivity (m)	1
Injection rate per well (m^3/d)	80
Pumping rate per well (m^3/d)	120
TCE injection concentration per well (mg/L)	50

2

3

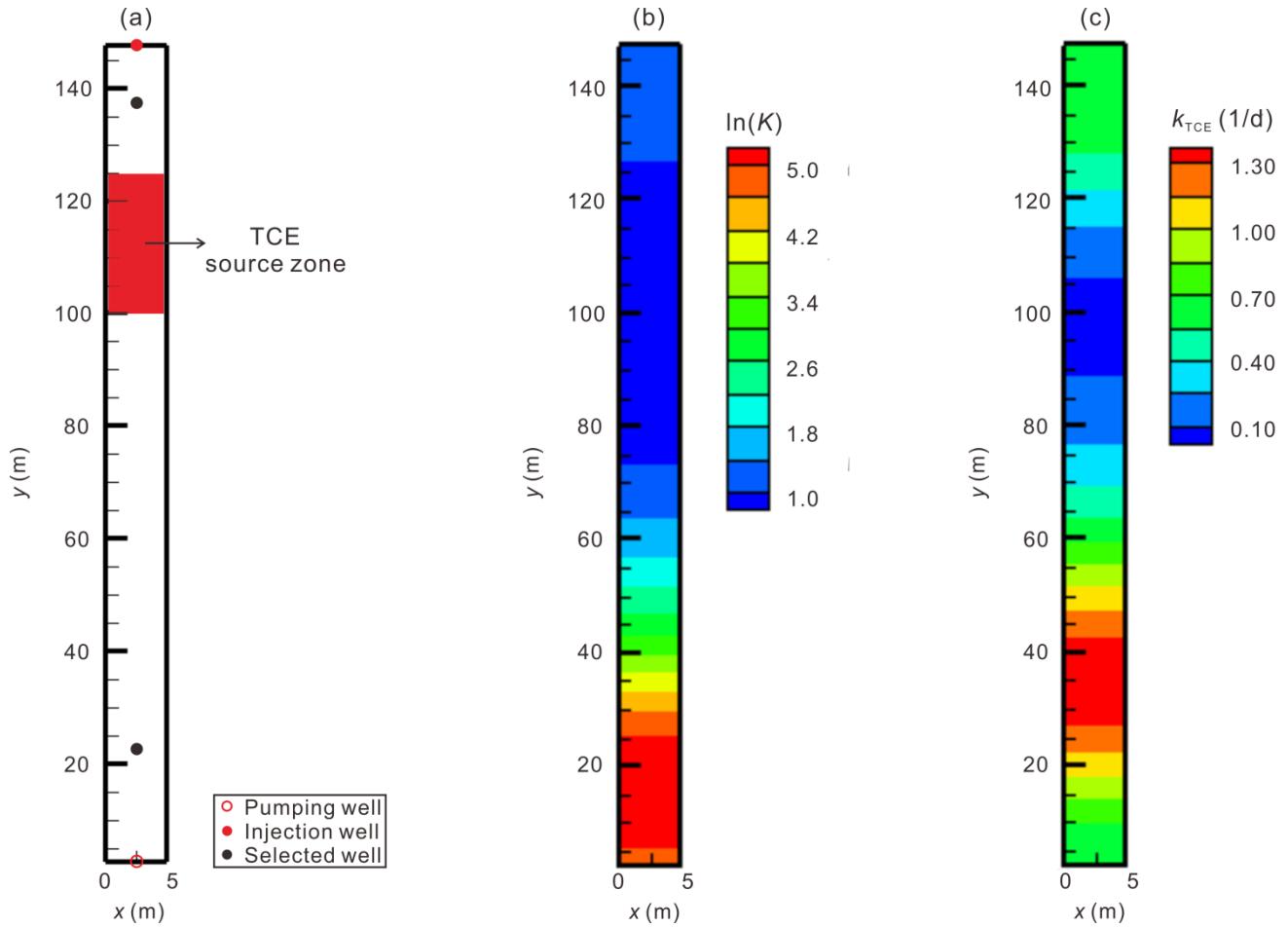
1

Table 4: Computational costs

Case name	Times of model invoking	Times of GA invoking	computing time (using the same computer)
Case 2	10	10	2 h
Case 14	8	2	26 min
Case 15	8	2	26 min
Case 16	12	3	25 min
Case 17	12	3	25 min

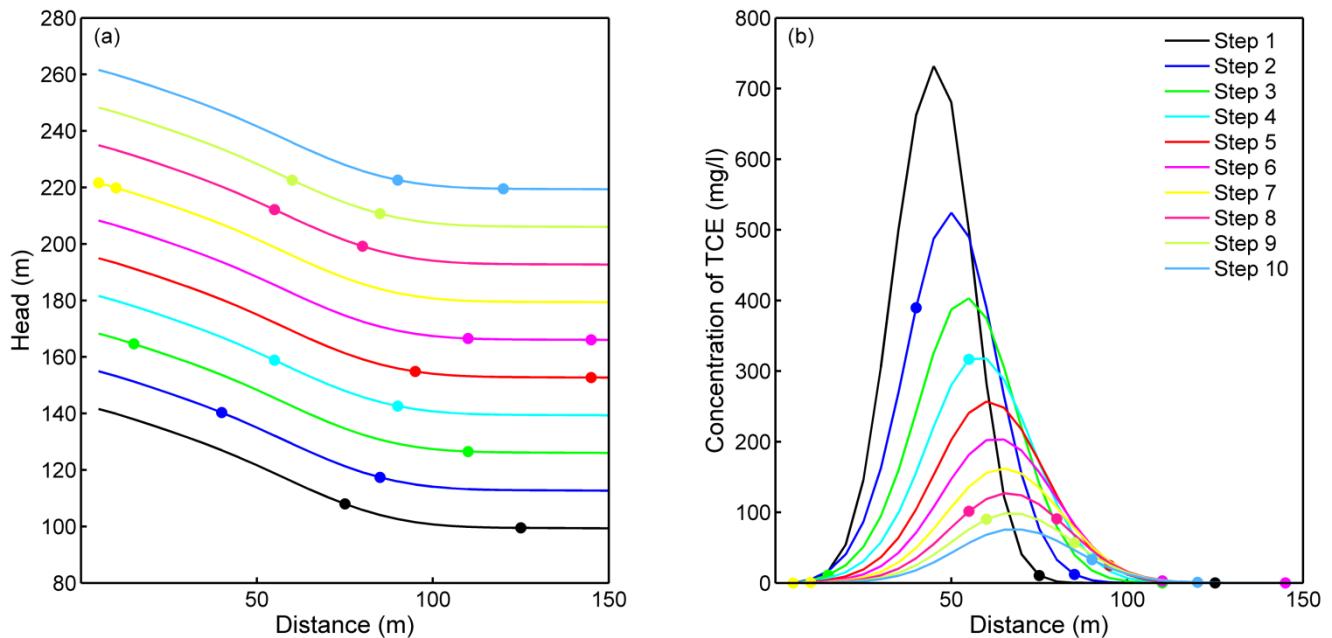
2

1 **Figures**

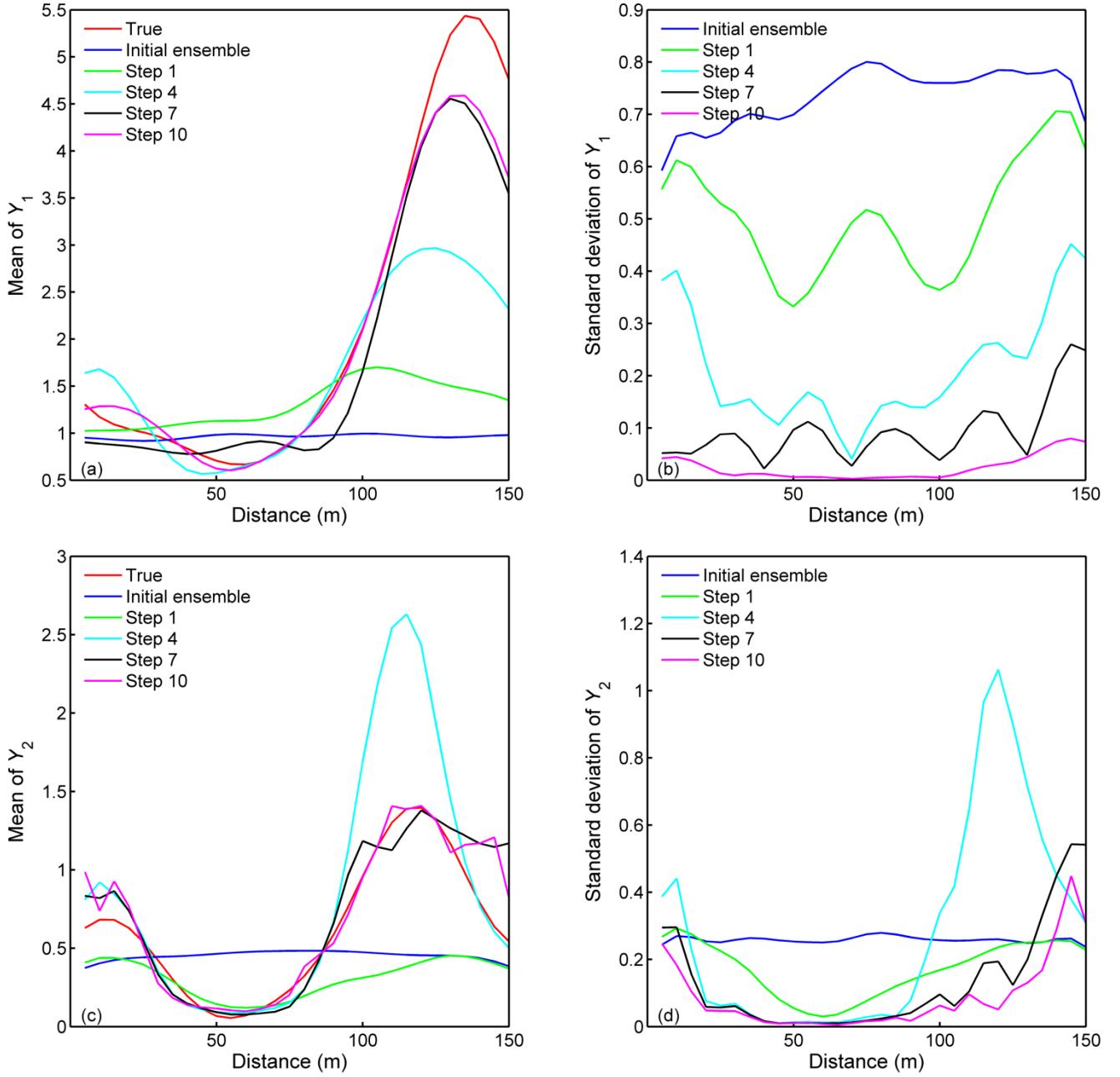


2
3 **Fig. 1** The conceptual model (a), the reference fields of the hydraulic conductivity (b) and the first-order
4 rate constant (c) for Case 1.

5



1 **Fig. 2** The calculated flow field (a) and concentration field (b) in Case 1. The circles denote the optimal
2 sampling locations proposed by the SEOD method.
3
4



1 **Fig. 3** The ensemble mean and the standard deviation of field Y_1 and Y_2 in Case 1. (a) and (b) are the
2 ensemble mean and the standard deviation of field Y_1 respectively, while (c) and (d) are the ensemble
3 mean and the standard deviation of field Y_2 respectively.
4

5

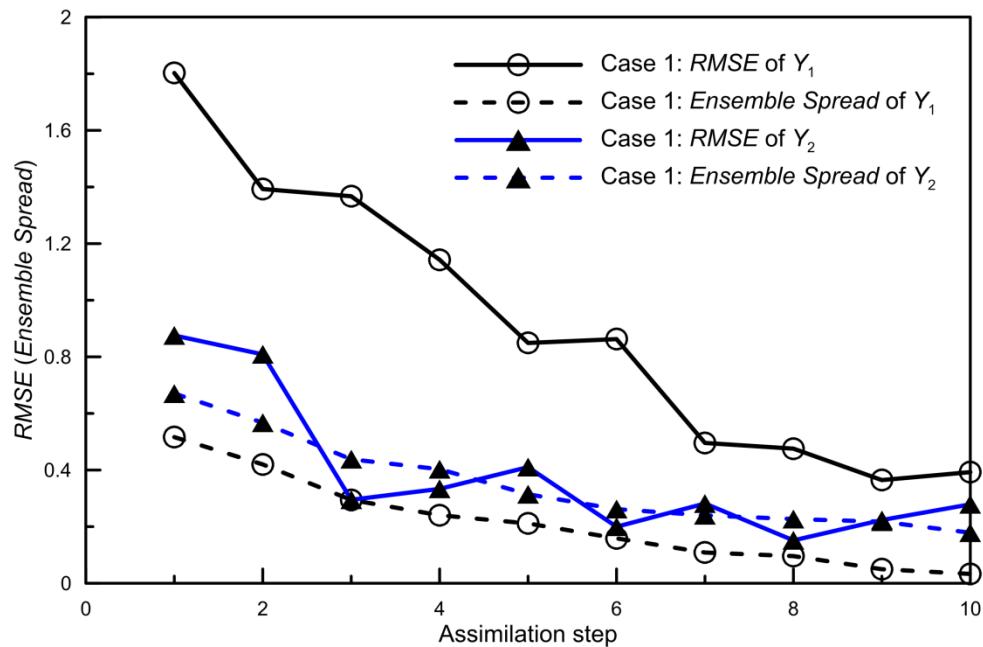


Fig. 4 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for field Y_1 and Y_2 in Case 1.

1
2
3
4

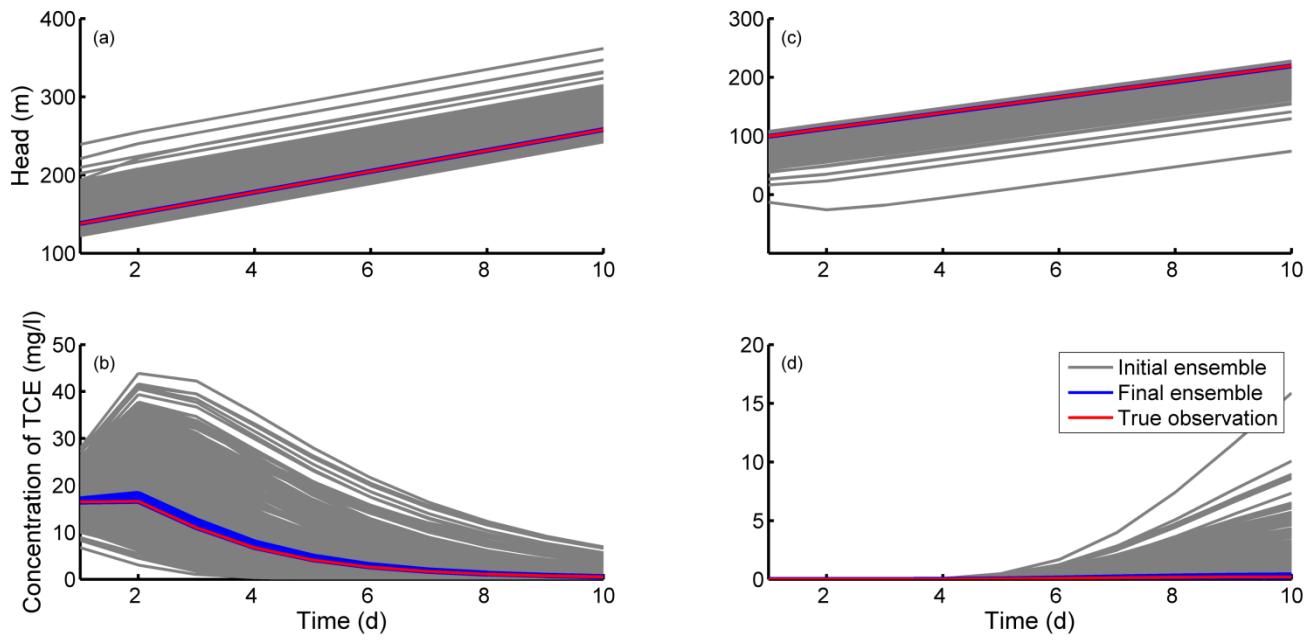
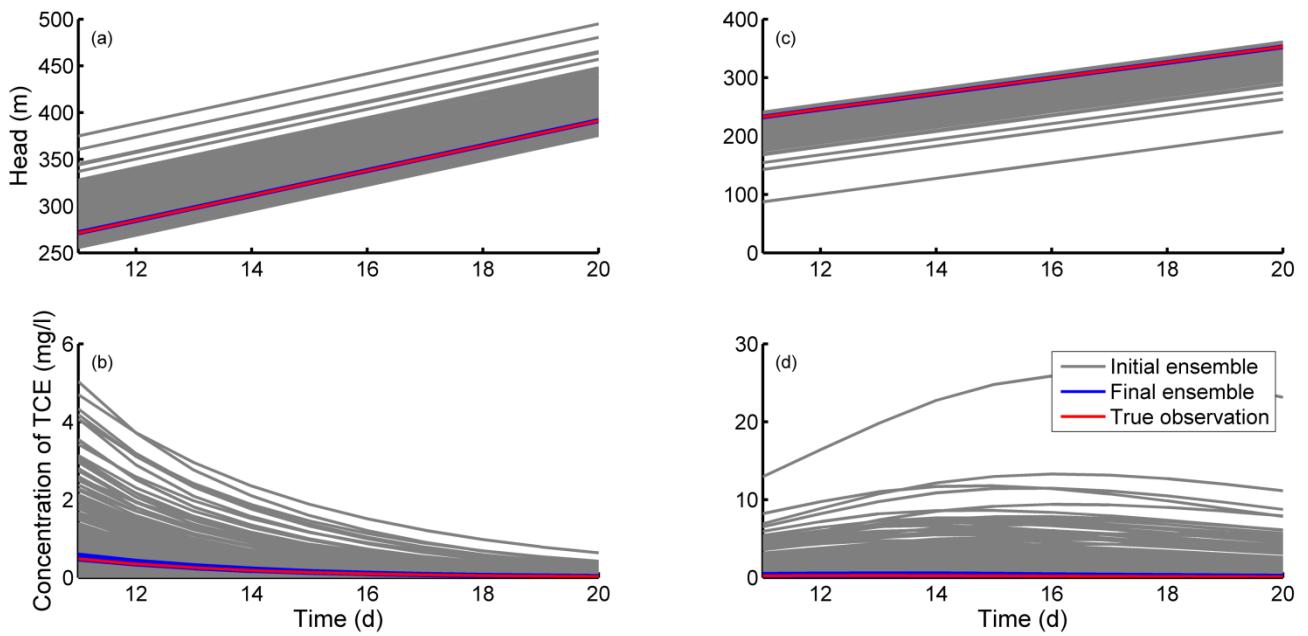


Fig. 5 The performance of data match. (a) and (c) show the data match of the head of two selected wells respectively, while (b) and (d) show the data match of the TCE concentration data of two selected wells respectively.



1 **Fig. 6** The performance of model prediction. (a) and (c) show the prediction of the head of two selected
2 wells respectively, while (b) and (d) show the prediction of the TCE concentration data of two selected
3 wells respectively.
4
5

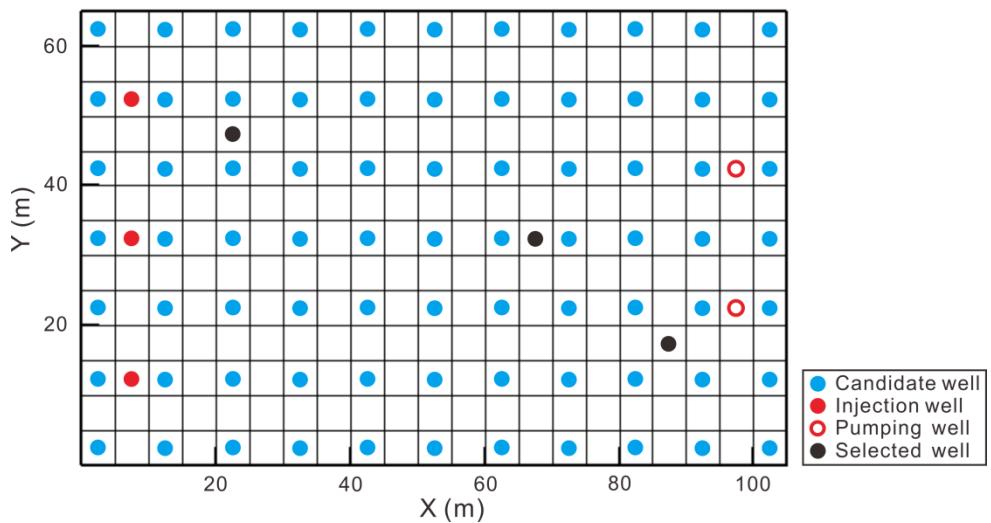
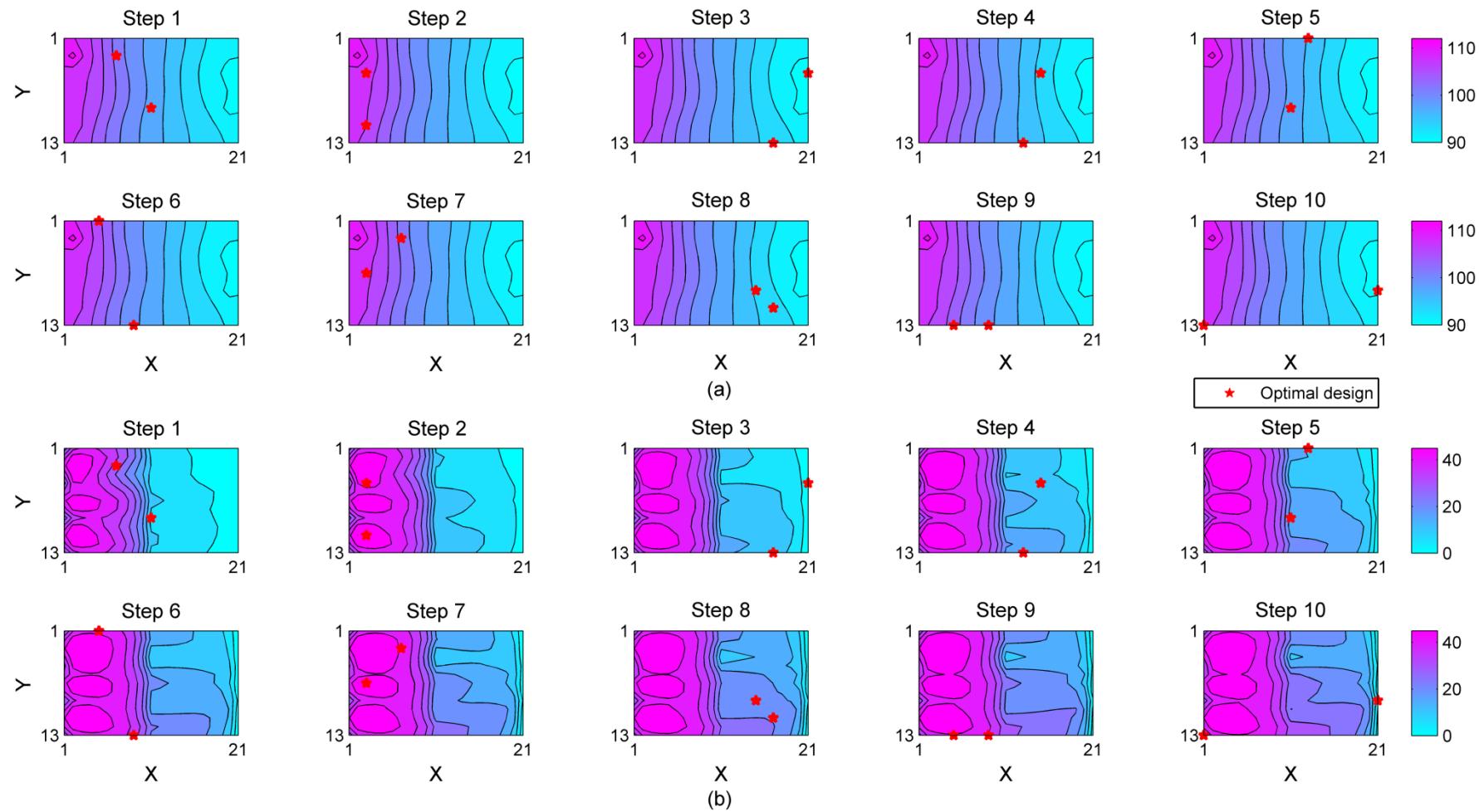


Fig. 7 The schematic of the two-dimensional conceptual model.

1
2
3



1
2 **Fig. 8** The optimal sampling locations (red stars) at every assimilation step of Case 2. (a) and (b) are the contour maps of the flow filed and the concentration
3 field, respectively.
4

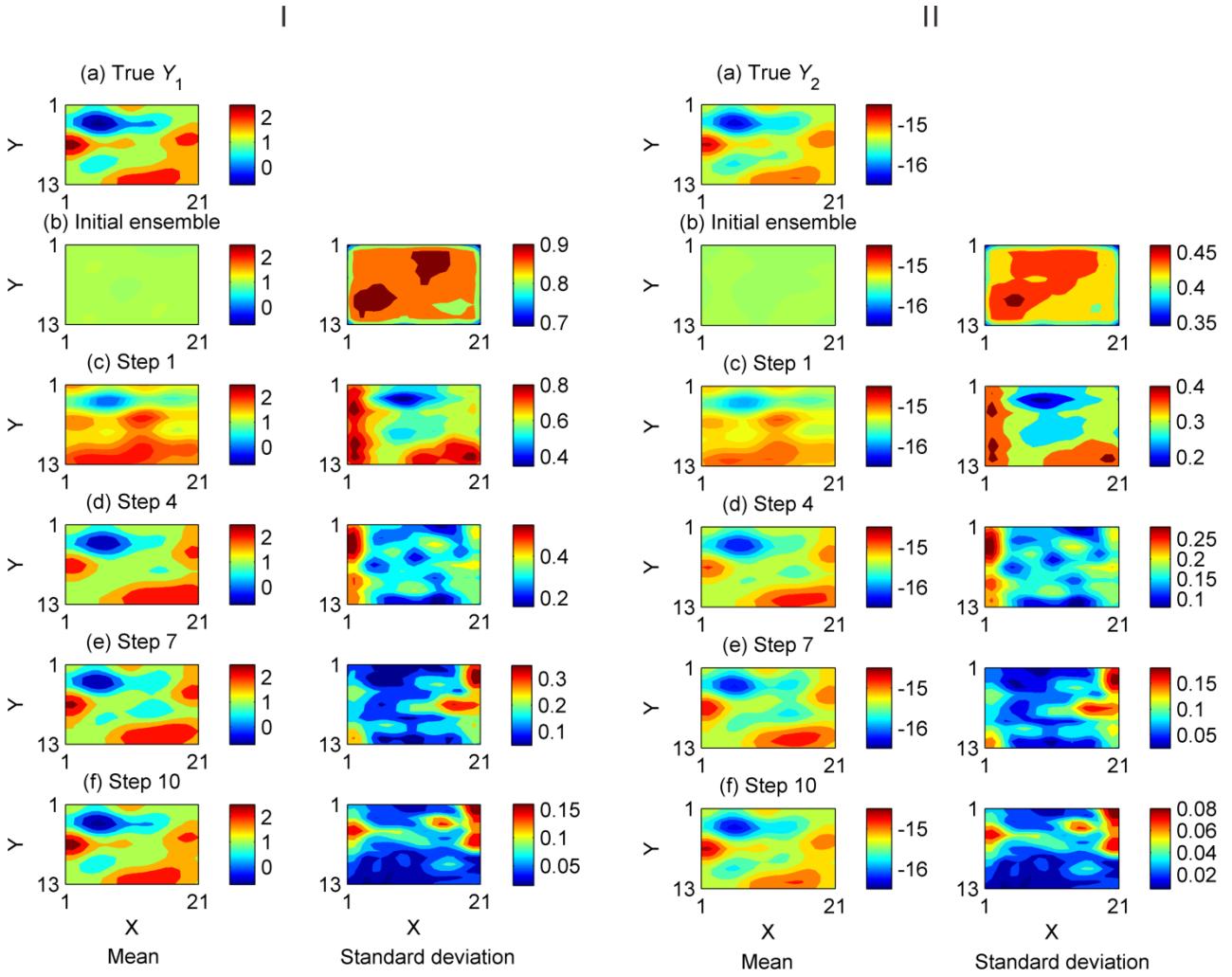


Fig. 9 The ensemble mean and the standard deviation of field Y_1 and Y_2 in Case 2. I for field Y_1 , II for field Y_2 .

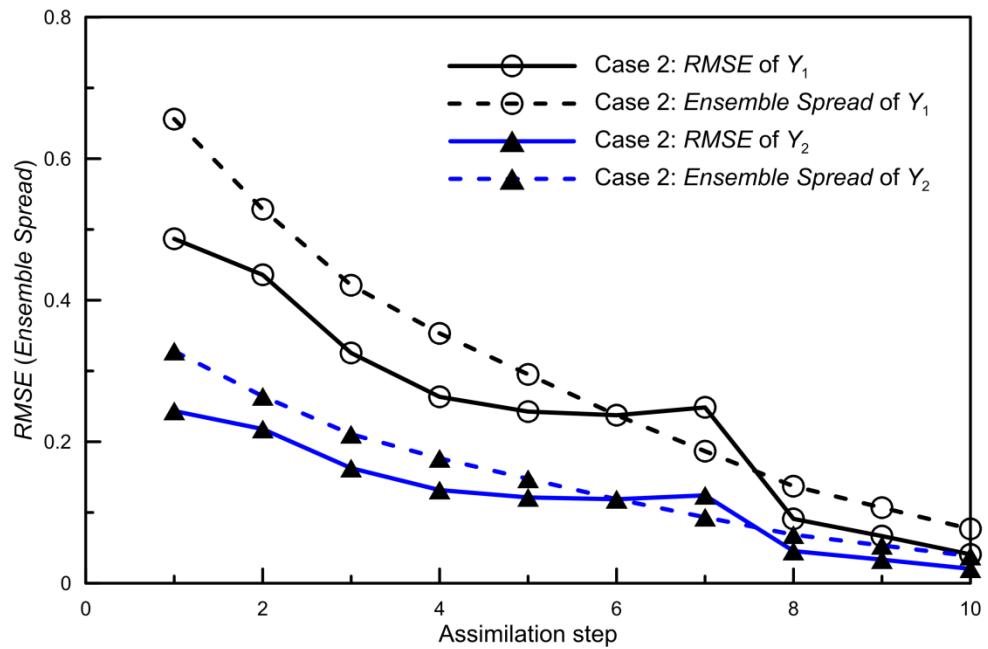


Fig. 10 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for field Y_1 and Y_2 in Case 2.

1

2

3

4

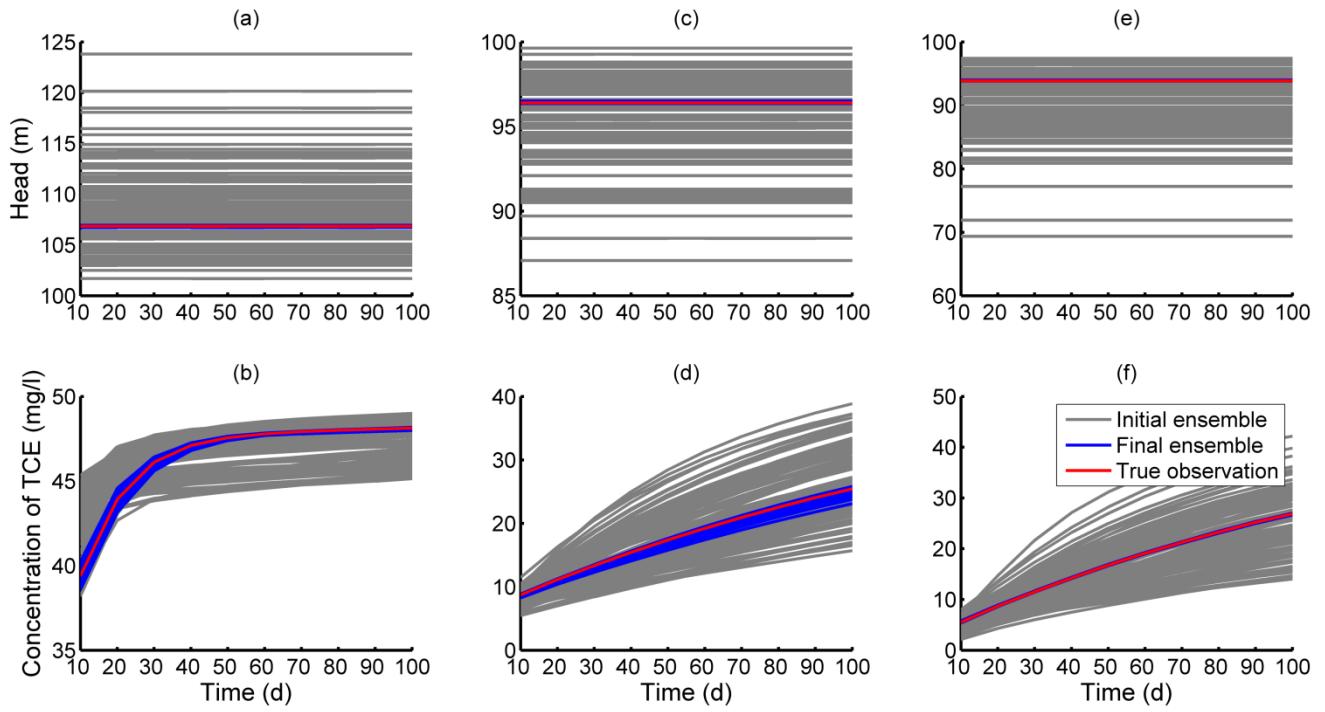


Fig. 11 The performance of data match. (a), (c) and (e) show the data match of the head of three selected wells respectively, while (b), (d) and (f) show the data match of the TCE concentration data of three selected wells respectively.

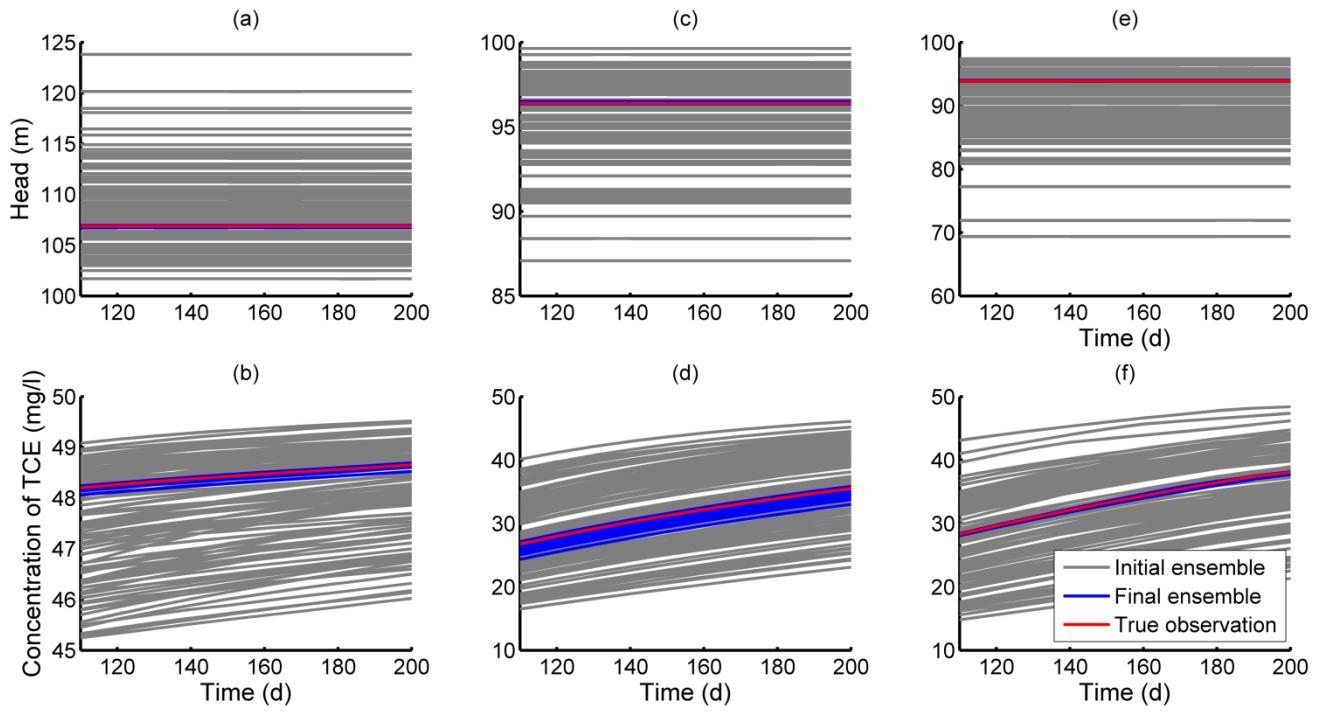


Fig. 12 The performance of model prediction. (a), (c) and (e) show the prediction of the head of three selected wells respectively, while (b), (d) and (f) show the prediction of the TCE concentration data of three selected wells respectively.

5

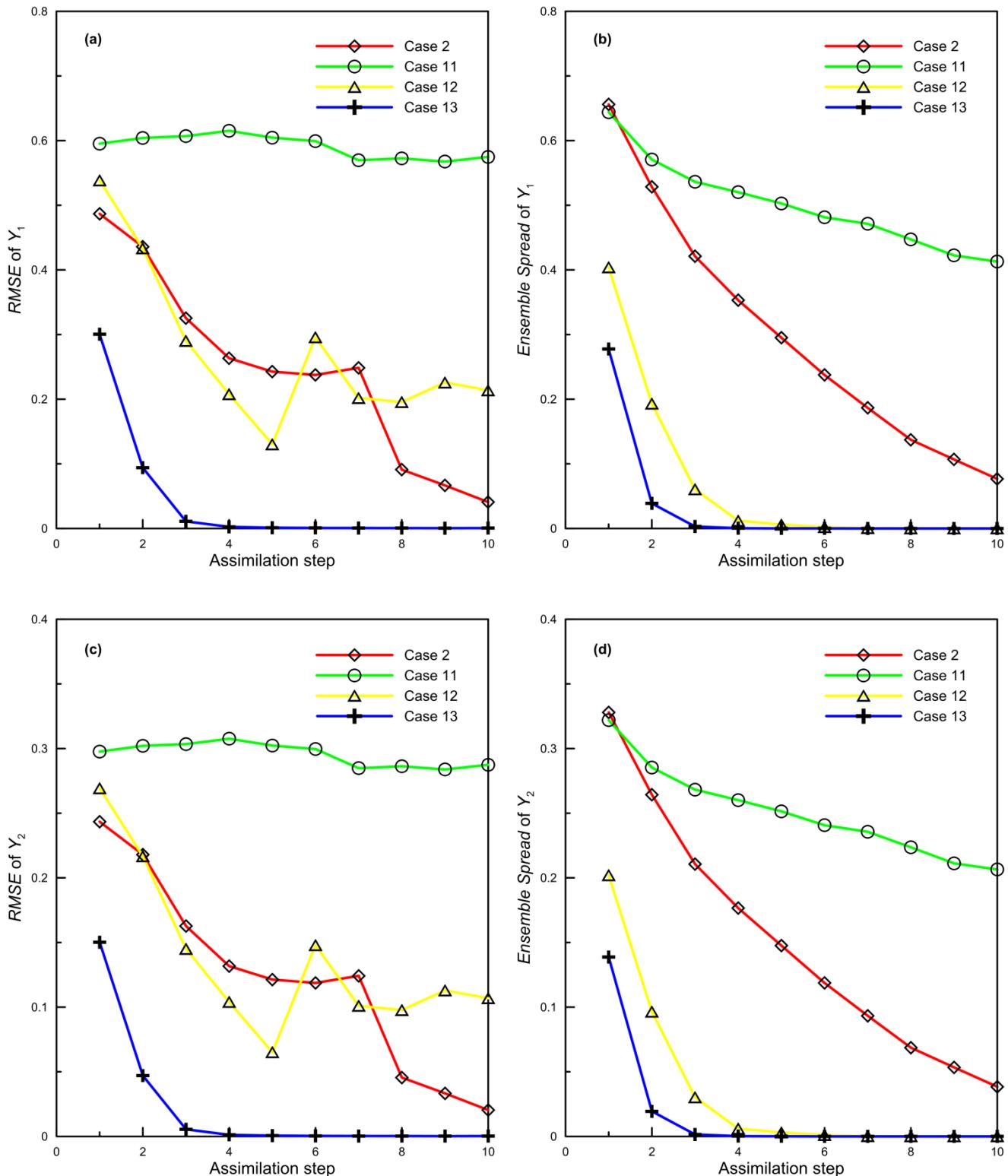


Fig. 13 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different sampling strategies.

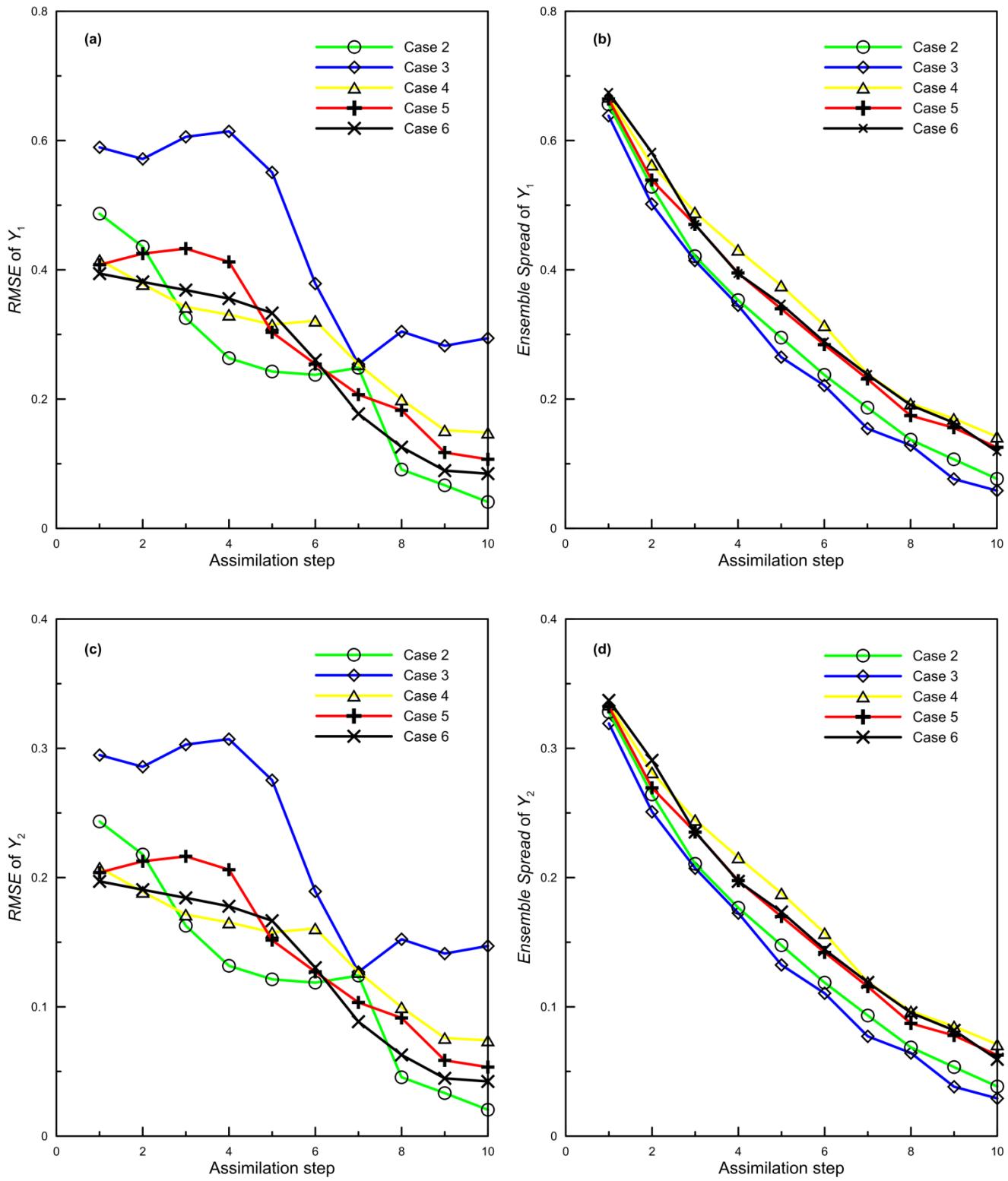
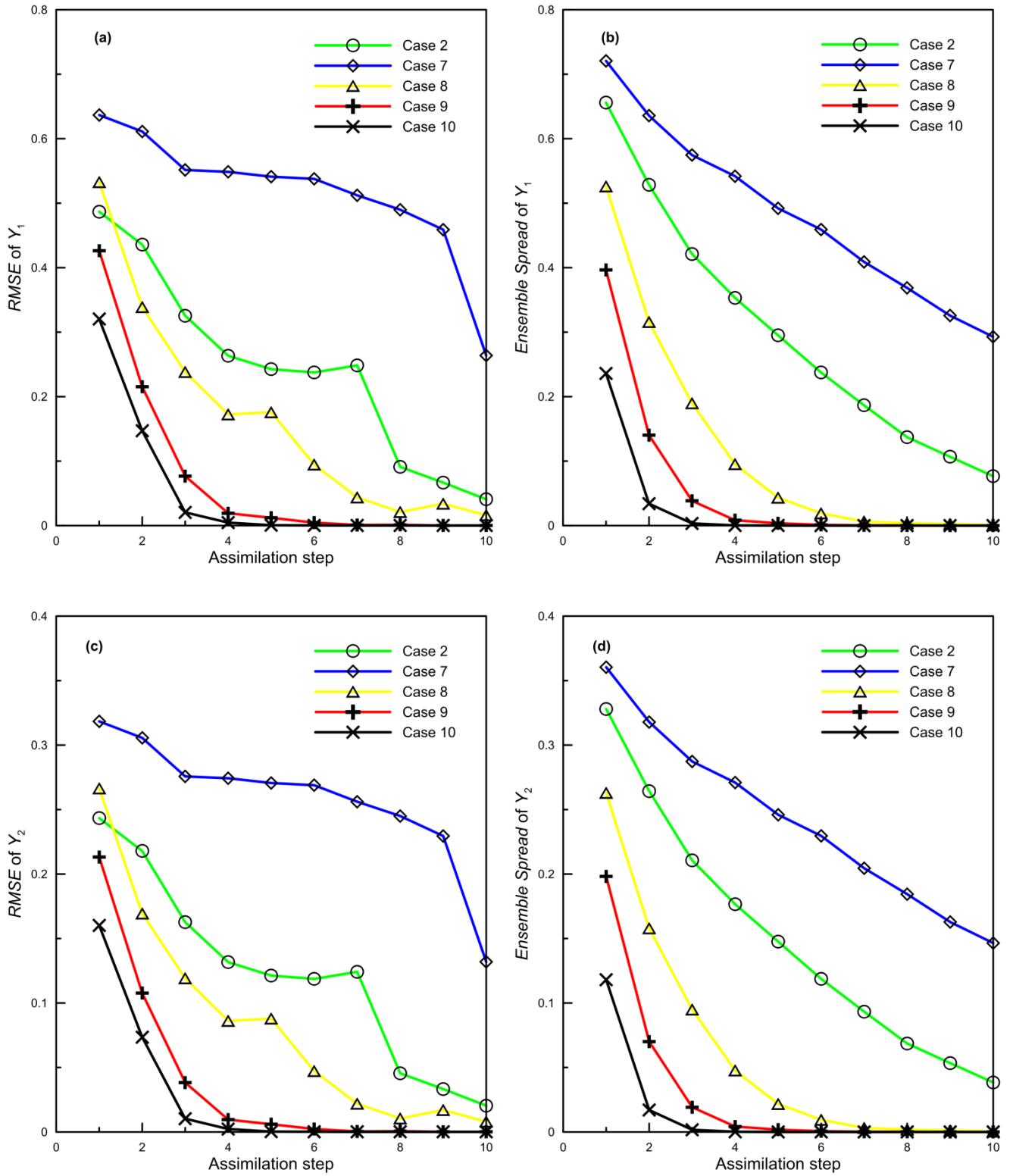


Fig. 14 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different ensemble sizes.



1 **Fig. 15** The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for
2 different numbers of optimal sampling locations.
3

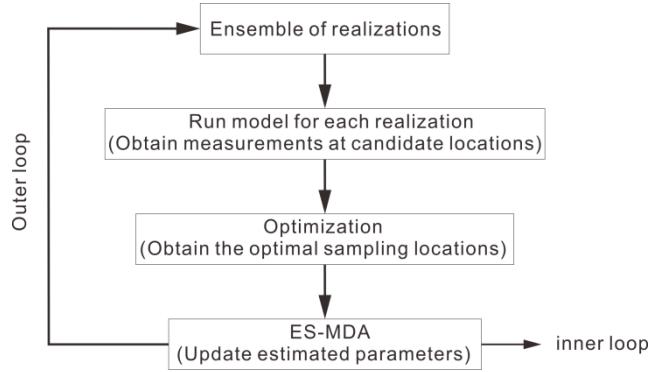


Fig. 16 The loop diagram of the improved SEOD method.

1

2

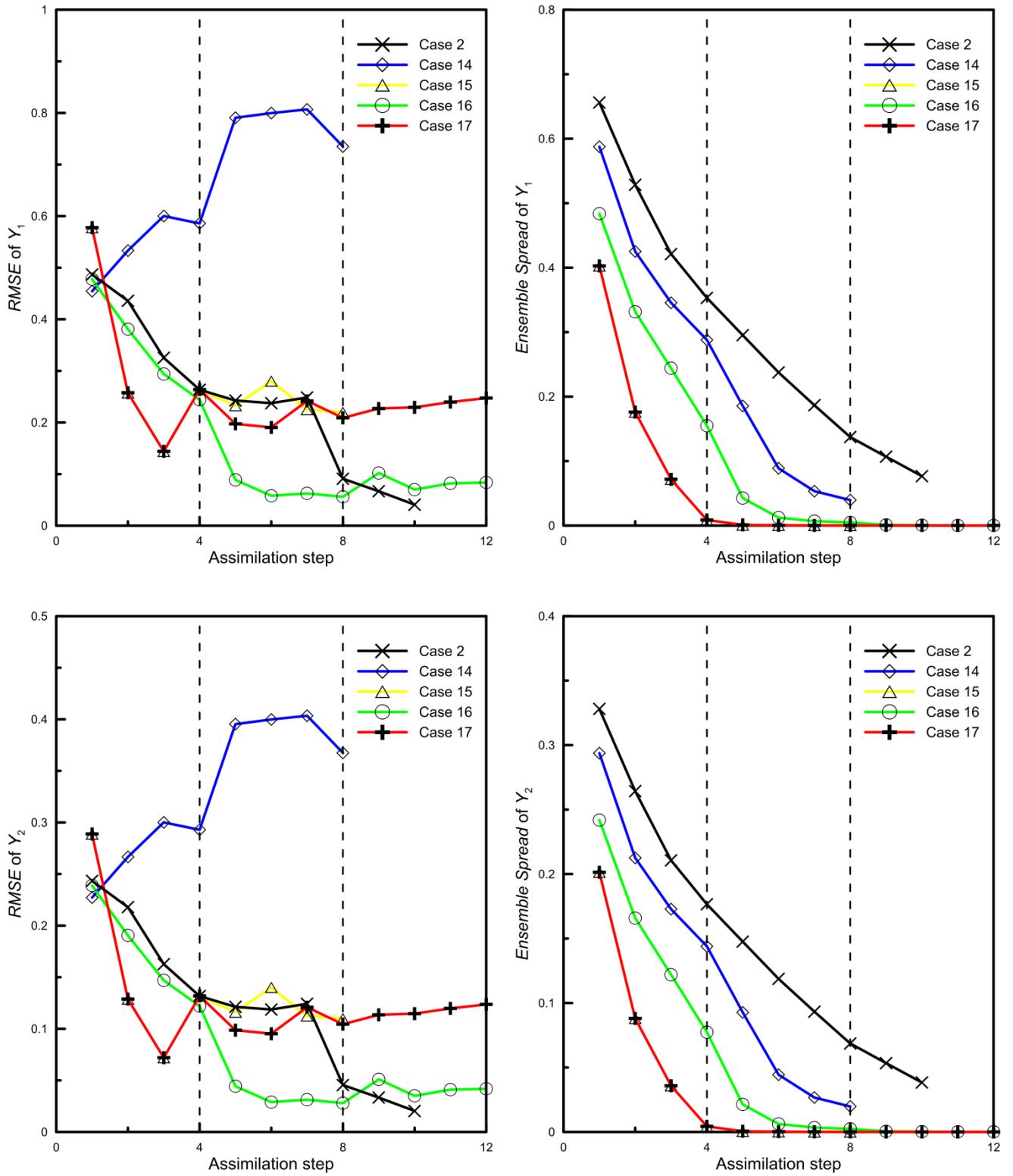


Fig. 17 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different cases using the improved SEOD method.



1 **Joint inversion of physical and geochemical parameters in
2 groundwater models by sequential ensemble-based optimal design**

3 Tian Lan¹, Xiaoqing Shi^{1,*}, Beilei Jiang², Yuanyuan Sun¹, Jichun Wu^{1,*}

4

5 ¹ Key Laboratory of Surficial Geochemistry, Ministry of Education and School of
6 Earth Sciences and Engineering, Nanjing University, Nanjing, China 210023

7 ² Nanjing Hydraulic Research Institute, National Key Laboratory of Water Resources
8 and Hydraulic Engineering, Nanjing, China, 210029

9

10 *Corresponding author:

11 Xiaoqing Shi, Email: shixq@nju.edu.cn; (86) 25-89680839

12 Jichun Wu, Email: jcwu@nju.edu.cn; (86) 25-89680705

13

14 **Highlights:**

15 1. The heterogeneous physical and geochemical parameters of groundwater models
16 are jointly estimated.

17 2. The SEOD method provides an efficient sampling strategy of sampling locations
18 to accurately estimate both physical and geochemical parameters.

19 3. The effectiveness and efficiency of the SEOD method for jointly estimating
20 physical and geochemical parameters are illustrated and demonstrated by the
21 comparison with EnKF using fixed sampling locations during estimation.

- 1 1. Both the heterogeneous physical and geochemical parameters of groundwater
- 2 models are accurately estimated by using the sequential ensemble-based optimal
- 3 design (SEOD) method.
- 4 2. The effectiveness and efficiency of the SEOD method are illustrated and
- 5 demonstrated by the comparison between the sequential optimization strategy
- 6 and the conventional strategy.
- 7 3. To enhance its computational efficiency, the SEOD method is improved by
- 8 replacing the EnKF with the Ensemble Smoother with multiple data assimilation
- 9 (ES-MDA).

10

1 **Abstract**

2 The Ensemble Kalman Filter (EnKF) method has been widely applied into the
3 parameter estimation of groundwater models. Its application in joint inversion of
4 physical and geochemical parameters of reactive transport models remains
5 challenging, partly due to the intrinsic heterogeneities of natural porous media and the
6 scarcity of observation data. Since it is generally cost intensive and time consuming to
7 acquire measurements from the subsurface system, it is beneficial to select the optimal
8 sampling locations that are most informative for model inversion. In this study, the
9 sequential ensemble based optimal design (SEOD) method is used to obtain the most
10 informative measurements at every assimilation step of EnKF. The performances of
11 SEOD are tested by jointly estimating heterogeneous physical and geochemical
12 parameters in two synthetic cases. The results show that the SEOD method provides
13 an effective designed sampling strategy of sampling locations to accurately estimate
14 heterogeneous distribution of physical and geochemical parameters. We also discuss
15 the impacts of ensemble size and number of optimal sampling locations on parameter
16 estimation. Meanwhile, the comparison between sequential optimization strategy and
17 conventional strategy (using fixed sampling locations during estimation) is also
18 investigated. Results show that larger ensemble size and more optimal sampling
19 locations improve the parameter estimation and the convergence of optimal sampling
20 strategy, however, the computational cost increases accordingly. Besides, the
21 effectiveness and efficiency of the SEOD method for jointly estimating physical and
22 geochemical parameters are also illustrated and demonstrated by the comparison with

1 EnKF using fixed sampling locations during estimation. Joint inversion of physical
2 and geochemical parameters in groundwater reactive transport models is still a great
3 challenge due to the intrinsic heterogeneities of natural porous media and the scarcity
4 of observation data. In this study, we make use of a sequential ensemble-based
5 optimal design (SEOD) method to jointly estimate physical and geochemical
6 parameters of groundwater models. The effectiveness and efficiency of the SEOD
7 method are illustrated by the comparison between the sequential optimization strategy
8 and the conventional strategy (using fixed sampling locations) for two synthetic cases.
9 Since the SEOD method is an optimization method based on the Ensemble Kalman
10 Filter (EnKF), it invokes the time-consuming Genetic Algorithm (GA) at every
11 assimilation step of the EnKF to obtain the optimal sampling locations. To enhance its
12 computational efficiency, we improve the SEOD method by replacing the EnKF with
13 the Ensemble Smoother with multiple data assimilation (ES-MDA). Furthermore, the
14 influence factors of the original and improved SEOD method are also discussed. Our
15 results show that the SEOD method provides an effective designed sampling strategy
16 to accurately estimate heterogeneous distribution of physical and geochemical
17 parameters. Moreover, the improved SEOD method is more advantageous than the
18 original one in computational efficiency, making this SEOD framework more
19 promising for future application.

20 **Keywords:** Optimal sampling strategy; Physical and geochemical heterogeneity;
21 Parameter estimation; Reactive transport model; Data assimilation.

1 **1 Introduction**

2 The subsurface environment is highly variable in its physical and chemical
3 composition. Heterogeneity of physical parameters (e.g. hydraulic conductivity) has
4 been shown to exert a key control on the mixing and spreading of conservative solutes
5 (Dagan 1984; Rubin 1991; Sudicky 1986). For reactive solutes, their transport and
6 reaction are simultaneously influenced by geochemical parameters (Atchley et al.
7 2014; Li et al. 2010; Scheibe et al. 2006). Like physical heterogeneity, the
8 heterogeneity of geochemical parameters exists as well, which could be caused, for
9 example, by spatial variability in the activity of bacteria related to biodegradation
10 (Fennell et al. 2001; Sandrin et al. 2004). Groundwater reactive transport model,
11 combining solute transport model and geochemical model, is an important tool to
12 characterize and explain the coupling processes. Identification of both physical and
13 geochemical parameters of the subsurface environment is critical for reliable
14 contaminant plume prediction, remediation and management. However, there is still
15 great challenge of jointly estimating the spatial distribution of physical and
16 geochemical parameters due to the intrinsic heterogeneities of natural porous media
17 and the scarcity of observation data.

18 Many ensemble based history matching methods have been developed and
19 widely applied in parameter estimation of groundwater models, such as Ensemble
20 Kalman Filter (EnKF) (Evensen 2009), ensemble smoother (Van Leeuwen and
21 Evensen 1996), and their iterative variants (Emerick and Reynolds 2013; Li and
22 Reynolds 2009). Among these ensemble based history matching methods, EnKF is

1 one of the most popular methods for its simplicity and flexibility in implementation.
2 EnKF is a variant of Kalman filter based on Monte Carlo method. In the EnKF
3 algorithm, it is not necessary to calculate the sensitivities explicitly which are required
4 in optimization-based approaches (Aanonsen et al. 2009, Zhou et al. 2014). Instead,
5 the cross correlations between model parameters and states calculated from
6 realizations are used (Evensen 2003, 2009). Besides, though EnKF is based on the
7 linear assumption, its efficiency and effectiveness in nonlinear problems with high
8 dimensionality have been illustrated (Chen and Zhang 2006; Hendricks Franssen and
9 Kinzelbach 2008; Moradkhani et al. 2005). To broaden its implementation scale and
10 improve its accuracy, many variants have been developed. For example, Gu and
11 Oliver (2007) proposed using the ensemble randomized maximum likelihood
12 (EnRML) method to obtain better results in strong non-linear models; Zhou et al.
13 (2011) proposed NS-EnKF method to solve non-Gaussian problems; Chen and Oliver
14 (2010), Emerick and Reynolds (2011) introduced a distance-based localization
15 scheme into EnKF to mitigate ensemble collapse.

16 Numerous variants based on EnKF are developed, however, few take sampling
17 design into consideration. In most of EnKF variants mentioned above, the sampling
18 locations were kept fixed during estimation, which is called conventional strategy in
19 the following discussion. It is intuitive that sampling locations predominantly
20 influence the data worth of measurements, and the parameter estimation can be
21 improved if the measurements are more informative even that the number of
22 measurements is the same. Knopman and Voss (1987, 1988) showed that

concentration data collected at locations that exhibited high sensitivities to model parameters provided the most information. Catania et al. (2004) built and tested a decision model which could obtain the "optimal" parameter estimation with minimum parameters uncertainty and experimental cost. Cleveland and Yeh (1990, 1991) presented two optimization algorithms in an aquifer tracer test for designing a monitoring network, which provided minimum parameter uncertainty. Though many studies have focused on designing optimal sampling strategy for minimizing parameter uncertainty in groundwater models (Carrera and Neuman 1986; Knopman and Voss 1989; Knopman et al. 1991; Nishikawa and Yeh 1989), few related researches have been studied within the EnKF framework. Man et al. (2016) developed a sequential ensemble based optimal design (SEOD) method, which firstly integrated sequential optimal design and information theory into the EnKF framework seamlessly to provide the most informative measurements for more accurate parameter estimation. He demonstrated the effectiveness of this method by estimating only physical parameters in unsaturated flow models, whereas we apply the SEOD method into jointly estimating physical and geochemical parameters of reactive transport models in this study.

In this work, the spatial distributions of physical and geochemical parameters are jointly estimated by assimilating both head and concentration measurements based on the SEOD method. The rest of the paper is organized as follows. In Section 2, the groundwater reactive transport model and the SEOD method are described. In Section 3, we construct synthetic one dimensional and two dimensional groundwater reactive

1 transport model cases. In Section 4, the effects of ensemble size and number of
2 optimal sampling locations are discussed. We also compare the sequential
3 optimization strategy and conventional strategy in section 4. Finally, conclusions are
4 summarized in Section 5. Joint inversion of physical and geochemical parameters in
5 groundwater reactive transport models is critical for reliable contaminant plume
6 prediction, remediation and management, but it is still a great challenge due to the
7 intrinsic heterogeneities of natural porous media and the scarcity of observation data.
8 The subsurface environment is highly variable in its physical and chemical
9 composition. Heterogeneity of physical parameters (e.g. hydraulic conductivity) has
10 been shown to exert a key control on the mixing and spreading of conservative solutes
11 (Dagan 1984; Rubin 1991; Sudicky 1986). For reactive solutes, their transport and
12 reactions are simultaneously influenced by geochemical parameters (Atchley et al.
13 2014; Li et al. 2010; Scheibe et al. 2006). Similar to physical heterogeneity, the
14 heterogeneity of geochemical parameters exists as well, which may be caused, for
15 example, by spatial variability in the activity of bacteria related to biodegradation
16 (Fennell et al. 2001; Sandrin et al. 2004). Therefore, it is important to jointly estimate
17 the spatial distribution of physical and geochemical parameters in groundwater
18 reactive transport models.

19 Inverse methods are often used by conditioning on observation data to
20 characterize the spatial variation of parameters, which has been extensively
21 investigated in the literature (e.g., Carrera et al. 2005; Dagan 1985; Doherty 2004;
22 Gómez-Hernández et al. 2003; Hendricks Franssen et al. 2009; Neuman 1973; Oliver

et al. 1997; Zhou et al. 2014). The Ensemble Kalman Filter (EnKF, Evensen 2003, 2009) is one of the most popular inverse methods over the last decade (Aanonsen et al. 2009; Oliver and Chen 2011; Zhou et al. 2014), recently used in parameter estimation and state prediction (Chen and Zhang 2006; Huang et al. 2009; Tong et al. 2010). It is a variant of the Kalman Filter (KF, Kalman 1960) based on the Monte Carlo method. Unlike the KF, the EnKF was developed for nonlinear problems (Evensen 2003, 2009), its efficiency and effectiveness in nonlinear problems with high dimensionality have been illustrated (Chen and Zhang 2006; Hendricks Franssen and Kinzelbach 2008; Moradkhani et al. 2005; Sorensen et al. 2004). In addition to the EnKF, the Ensemble Smoother (ES, Van Leeuwen and Evensen 1996) and its iterative variants, like the Ensemble Smoother with multiple data assimilation (ES-MDA, Emerick and Reynolds 2013), are popular as well. Unlike the EnKF, the ES and the ES-MDA perform global update rather than sequential update during the data assimilation, avoiding restarting models again and again, so they are of more simplicity and computational efficiency than the EnKF.

Much research has focused on developing better methods based on the EnKF to broaden its implementation scale and improve its accuracy (Chen and Oliver 2010; Emerick and Reynolds 2011; Gu and Oliver 2007; Li and Reynolds 2009), with the sampling locations fixed during the data assimilation (called the conventional strategy in the following discussion). However, it is intuitive that the data worth of measurements is dramatically influenced by sampling locations, and the parameter estimation result can be improved if the measurements are more informative even

1 though the number of sampling locations is the same. There has been much research
2 revealed the effect of sampling strategies on the parameter uncertainty and predictive
3 uncertainty in groundwater models (Carrera and Neuman 1986; Cleveland and Yeh
4 1990; Knopman and Voss 1987; Nowak et al. 2010; Sun and Yeh 2007; Ushijima and
5 Yeh 2015; Zhang et al. 2015). In view of these two aspects, Man et al. (2016)
6 integrated a sequential optimal design and the information theory into the EnKF
7 framework seamlessly to provide the most informative measurements for more
8 accurate parameter estimation, and proposed a sequential ensemble-based optimal
9 design (SEOD) method. Man et al. (2016) demonstrated the effectiveness of this
10 method by estimating only physical parameters in unsaturated flow models,
11 assimilating only piezometric head data. However, the SEOD method developed by
12 Man et al. (2016) invokes the optimization algorithm (the Genetic Algorithm) at each
13 assimilation step, so its computational efficiency is not very satisfying. Furthermore,
14 to the best of our knowledge, few studies have focused on joint inversion of physical
15 and geochemical parameters by assimilating multiple kinds of data.

16 The objective of this study is to estimate both physical and geochemical
17 parameters accurately in groundwater models by using the recent proposed SEOD
18 method, and to enhance the computational efficiency of the SEOD method by
19 replacing the EnKF with the ES-MDA. The rest of the paper is organized as follows.
20 In Section 2, the groundwater reactive transport model and the SEOD method are
21 described. In Section 3, synthetic one-dimensional and two-dimensional groundwater
22 reactive transport model cases are constructed to jointly estimate the physical and

1 geochemical parameters by using the SEOD method. In Section 4, the comparison
 2 between the sequential optimization strategy and the conventional strategy, and the
 3 effects of the ensemble size and the number of optimal sampling locations are
 4 discussed. Furthermore, we improve the SEOD method by replacing the EnKF with
 5 the ES-MDA to enhance its computational efficiency, and make a comparison of the
 6 original and the improved SEOD method in Section 4.4. Conclusions are summarized
 7 in Section 5.

8 **2 Methodologies**

9 **2.1 Groundwater reactive transport model**

10 In this work, transient flow is assumed, as the following governing equation
 11 (Bear 1972),

$$12 \quad \nabla \cdot (K \nabla H) + W = \mu_s \frac{\partial H}{\partial t} \quad (1)$$

13 where $\nabla \cdot$ is the divergence operator; ∇ is the gradient operator; K is the hydraulic
 14 conductivity [LT^{-1}]; H is the hydraulic head [L]; W is the volumetric injection
 15 (pumping) flow rate per unit volume of the aquifer [LT^{-1}]; μ_s is the specific storage
 16 coefficient of the aquifer [L^{-1}]; t is the time [T].

17 The governing equation for the transport and reactions of aqueous species is
 18 defined as (Zheng 2006; Prommer and Post 2010):

$$19 \quad \frac{\partial C_n}{\partial t} = \nabla \cdot (D \cdot \nabla C_n) - \nabla \cdot (v C_n) + r_{reac,n} + \frac{q_s}{\theta} C_n^s \quad (2)$$

20 where C_n is the aqueous concentration of the n th component [ML^{-3}]; t is the time [T];
 21 D is the diffusion coefficient [L^2T^{-1}]; $v = (-K \nabla H) / \theta$ [L^2T^{-1}]; $r_{reac,n}$ is the
 22 concentration change of the n th component caused by reactions; q_s is the volumetric

1 flow rate per unit volume of the aquifer [T^{-1}]; θ is the effective porosity; and C_n^s is the
2 concentration of the source or sink flux of the n th component [ML^{-3}].

3 Eq. (1) is solved by the numerical code MODFLOW-2000 (Harbaugh et al.
4 2000), and Eq. (2) is solved by the numerical code MT3DMS (Zheng 2006).

5 **2.2 Sequential ensemble-based optimal design (SEOD) method**

6 The sequential ensemble-based optimal design (SEOD) method is a new **recently**
7 **proposed** optimal method based on **the** EnKF, **proposed recently** (Man et al. 2016). At
8 each recursive step, the SEOD method provides an optimal sampling strategy, giving
9 the maximum value of information metric. Then, the analysis equation of **the** EnKF is
10 used to update estimated parameters by assimilating the most informative
11 measurements, obtained based on the optimal sampling strategy.

12 In this work, relative entropy (RE), also known as the Kullback-Leibler
13 divergence (Kullback 1977), is used to measure the information content of the
14 posterior probability density function (pdf) relative to the prior pdf. If these two
15 distributions are both n -dimensional Gaussian, RE between these two distributions is
16 defined as **follows**:

$$RE = J_b + [\ln \det(\mathbf{B}\mathbf{A}^{-1}) + \text{Tr}(\mathbf{A}\mathbf{B}^{-1}) - n] / 2 \quad (3)$$

17 where $J_b = (\mathbf{a} - \mathbf{b})^T \mathbf{B}^{-1} (\mathbf{a} - \mathbf{b}) / 2$ is the signal part of RE ; $\det(\cdot)$ denotes the
18 **determinant**, $\text{Tr}(\cdot)$ denotes the trace; \mathbf{a} and \mathbf{A} denote the mean and covariance matrix
19 of prior statistics respectively; \mathbf{b} and \mathbf{B} denote the mean and covariance matrix of
20 posterior statistics respectively.

22 **The loop of the SEOD method for parameter estimation in details can be found**

1 in Man et al. (2016). The loop of the SEOD method for parameter estimation is briefly
 2 recalled. More details can be found in Man et al. (2016). In the EnKF, all the
 3 parameters of interest \mathbf{p} are augmented with state variables \mathbf{h} into a joint state vector
 4 $\mathbf{x} = [\mathbf{p} \; \mathbf{h}]^T$. Before the forecast step, an ensemble of N_e realizations of parameters is
 5 generated.

(I) Forecast step

7 Rerun the forward model G from time 0 to time step $j+1$ with parameters updated
8 at time step j (Eq. (4)).

$$\mathbf{x}_{i_{j+1}}^f = G(\mathbf{x}_{i_j}^a), \quad i=1,2,\dots,N_e \quad (4)$$

In the above equation, i is the ensemble member index, j is the time step index, superscripts f and a denote forecast and analysis, respectively.

(II) Optimal design

Given a specific sampling strategy \mathbf{H}' , the possible realizations of measurements can then be expressed as $\mathbf{d}'_i = \mathbf{H}'\mathbf{x}_i^f + \xi_i$. With the realizations of measurements, the updated ensemble can be obtained from the Eq. (6). According to the prior and posterior statistics (mean and covariance), the information metrics RE of each candidate sampling strategies can be calculated. By comparing the RE values of different candidate sampling strategies, the optimal sampling design \mathbf{H}_{opt} can be determined by solving the following optimization problem (Eq. (5)) with the help of the Genetic Algorithm (GA, Whitley 1994).

$$21 \quad \mathbf{H}_{\text{opt}} = \arg \max RE(\mathbf{H}) \quad (5)$$

22 (III) Analysis step

1 After obtaining the optimal sampling strategy, the actual measurements \mathbf{d} can
 2 be obtained and used in the analysis step (Eq. (6)).

3
$$\mathbf{x}_{i,j+1}^a = \mathbf{x}_{i,j+1}^f + \mathbf{C}_{YD}(\mathbf{C}_{DD} + \mathbf{C}_D)^{-1}(\mathbf{d}_{obs_{i,j+1}} - \mathbf{d}_{i,j+1}), \quad i=1,2,\dots,N_e \quad (6)$$

4 In the above equation, \mathbf{C}_{YD} is the cross-covariance matrix between the forecast
 5 state and the predicted data, \mathbf{C}_{DD} is the covariance matrix of the predicted data, \mathbf{C}_D
 6 is the covariance matrix of the measurements error, \mathbf{d}_{obs} is the perturbed
 7 observations with noise of covariance \mathbf{C}_D , and \mathbf{d} is the predicted data.

8 After the analysis step, the updated ensemble \mathbf{X}^a is obtained. Then, go back to
 9 step (I), the updated ensemble obtained this step is implemented for the next step.

10 To evaluate the performance of parameter estimation, two commonly used
 11 indicators, the *RMSE* and the *Ensemble Spread*, are defined as follows:

12
$$RMSE = \sqrt{\frac{1}{N_m} \sum_{i=1}^{N_e} (\bar{Y}_i - Y_i)^2} \quad (74)$$

13
$$Ensemble \quad Spread (ES) = \sqrt{\frac{1}{N_m} \sum_{i=1}^{N_e} \text{var}(Y_i)} \quad (85)$$

14 where \bar{Y} and Y are the estimated and the reference field respectively; $\text{var}(Y)$ is
 15 the ensemble variance of the field; N_m is the total number of nodes in the study
 16 domain; N_e is the ensemble size; i is the node index. The *RMSE* measures the
 17 accuracy of the estimation, while the *Ensemble Spread* measures the uncertainty
 18 in of the estimation.

19 **3 Case studies**

20 **3.1 Case 1: One-dimensional synthetic case**

21 In this case, a one-dimensional confined aquifer with a starting head of 100 m

1 is constructed, in which saturated transient flow is assumed. As shown in Fig.
 2 1(a), we choose the horizontal aquifer is to be 5 m×150 m and the grid space is to be 5
 3 m both in horizontal x and y direction. Then, a Trichloroethylene (TCE) leaking
 4 area with an initial concentration of 1000 mg/L is introduced into the aquifer, and the degradation of TCE is regarded as assumed to follow the
 5 first-order kinetic reaction. Furthermore, an injection well and a pumping well are
 6 set upstream and downstream respectively, and all boundaries of the aquifer are
 7 assumed to be impermeable. In this case, the spatial distribution of the hydraulic
 8 conductivity (K) and the first-order rate constant (k_{TCE}) (Fig. 1(b), (c)) are jointly
 9 estimated. At every assimilation step, 2 optimal sampling locations, are selected from
 10 30 candidate locations, to provide the most informative measurements.
 11 More details are given in Table 1 and 2.

13 [Figure 1]

14 [Table 1]

15 [Table 2]

16 The log saturated hydraulic conductivity $Y_1 = \ln(K)$ and the first-order rate
 17 constant $Y_2 = k_{TCE}$ are assumed to be Gaussian distributed, with mean $\mu_{Y_1}=1$ and
 18 $\mu_{Y_2}=0.17$ and variance $\sigma_{Y_1}=1$ and $\sigma_{Y_2}=0.47$ respectively. Two arbitrary locations
 19 (x_1, y_1) and (x_2, y_2) in the random field are assumed to be correlated in the following
 20 form:

$$21 C_Y(x_1, x_2) = C_Y(x_1, y_1; x_2, y_2) = \sigma^2 \exp\left[-\frac{|x_1 - x_2|}{\lambda_x} - \frac{|y_1 - y_2|}{\lambda_y}\right] \quad (9)$$

22 where the horizontal correlation length $\lambda_x= 5$ m, and the vertical correlation length λ_y

1 = 30 m. Here we use the Karhunen-Loeve (K-L) expansion (Zhang and Lu 2004) to
2 parameterize the random field so as to achieve the reference fields and initial
3 ensemble members. The measurement errors of **the** head and concentration **data** are
4 assumed to follow **the** standard normal distribution with the standard deviation of 0.01
5 m and 10^{-6} mg/L respectively. Since the SEOD method is sequential, the uncertainty
6 changes as real-time measurements are assimilated, which leads to **the** optimal
7 sampling locations changing with time. The optimal sampling locations at 10
8 assimilation steps are shown in Fig. 2. It shows that **the** optimal sampling locations
9 change with the flow and concentration field so as to obtain the most informative
10 measurements. It is **of** interesting to note that most optimal sampling locations are
11 **around** located at the front of **the** contaminant plume, which means that **sampling**
12 **locations around the contaminant plume front** **these locations** can provide the most
13 informative measurements.

14 [Figure 2]

15 In Fig. 3, we plot the curves of **the** ensemble mean and **the** standard deviation at
16 different assimilation-**stages** **steps**. It shows that, for both Y_1 and Y_2 , the ensemble
17 mean at **the** final assimilation step is very close to the reference field. **In addition**
18 **Furthermore**, the ensemble standard deviation is high at **the** early steps; however, it
19 reduces dramatically after assimilating the most informative measurements from **the**
20 optimal sampling locations. As shown in Fig. 4, the *RMSE* and **the Ensemble Spread**
21 of Y_1 and Y_2 decrease as **the** assimilation step increases, which also suggests that
22 estimated fields are close to their reference fields and of low uncertainty.

1 [Figure 3]

2 [Figure 4]

3 To further illustrate the accuracy and uncertainty of the estimation, we also
4 evaluate the performance of data match and model prediction. Considering the
5 limitation of space, two wells, (marked with black circles in Fig. 1(a)), are selected
6 randomly to show the following evaluation results.

7 The initial and final ensembles of Y_1 and Y_2 are taken into the synthetic model
8 respectively to calculate the head and concentration data, which are then compared
9 with the real observations. After data assimilation, the data calculated by the final
10 ensemble become closer to the real observations, as shown in Fig. 5 and 6. Overall,
11 model uncertainty is significantly reduced after assimilating the most informative
12 measurements from the optimal sampling locations.

13 [Figure 5]

14 [Figure 6]

15 **3.2 Case 2: Two-dimensional synthetic case**

16 In this case, saturated transient flow is assumed in the a two-dimensional
17 confined aquifer with the a starting head of 50 m. As shown in Fig. 7, we choose the
18 horizontal aquifer is to be 105 m×65 m and the grid space is to be 5 m both in
19 horizontal x and y direction. Three TCE leaking sources with the constant injection
20 flow of 80 m³/d and the constant concentration of 50 mg/L per well are set upstream
21 in the aquifer, and Furthermore, the liner sorption reaction of the TCE is considered
22 in this case. Besides three injection wells, two pumping wells with the constant

1 pumping flow of 120 m³/d per well are set downstream (Fig.7). All In addition, all
2 boundaries of the aquifer are assumed to be impermeable. In this case, the hydraulic
3 conductivity (K) and the liner sorption constant (k_d), parameters to be estimated, are
4 assumed to be spatially heterogeneous (Fig.9). At every assimilation step, 2 optimal
5 sampling locations are selected from 77 candidate locations (Fig.7) to provide the
6 most informative measurements. More details are given in Table 2 and 3.

7 [Figure 7]

8 [Table 3]

9 The log saturated hydraulic conductivity $Y_1 = \ln(K)$ is assumed to be Gaussian
10 distributed with mean $\mu_{Y_1}=1$ and variance $\sigma_{Y_1}=1$. The horizontal and vertical
11 correlation lengths of Y_1 are 40 m and 20 m respectively. With these statistics, the
12 reference field and initial ensemble members of Y_1 can be generated by the K-L
13 (Karhunen-Loeve) decomposition based on Eq. (96). For field Y_2 , it is assumed that
14 there is a positive correlation between $Y_2 = \ln(k_d)$ and Y_1 , i.e. $Y_2 = 0.5 \times Y_1 - 15.95$, on
15 which the generation of reference field Y_2 and its initial ensemble members are based
16 on. In addition, the measurement errors of the head and concentration data are
17 assumed to follow the standard normal distribution with the standard deviation of 0.01
18 m and 10⁻⁶ mg/L respectively.

19 In Fig. 8, we plot the optimal sampling locations at each assimilation step. It
20 shows the tendency that locations of large gradient are more likely to be selected as
21 the optimal sampling locations. Overall, it shows that the choice of the optimal
22 sampling locations at each assimilation step changes with the flow and concentration

1 field so as to obtain the most informative measurements.

2 [Figure 8]

3 The contour maps of the ensemble mean and the standard deviation at different
4 assimilation stages steps are plotted in Fig. 9. It is shown shows that, for both Y_1 and
5 Y_2 , the contour maps of the ensemble mean exhibit a pattern very similar to the
6 reference fields. Even just after 4 assimilation steps, the contour maps of the ensemble
7 mean recover the major features of the reference fields of Y_1 and Y_2 . Furthermore
8 Besides, the ensemble standard deviations reduce dramatically after assimilating the
9 most informative measurements from the optimal sampling locations, indicating that
10 the optimal sampling strategy does play it's a crucial role in the model inversion
11 though only 2 sampling locations are selected at every assimilation step.

12 [Figure 9]

13 The RMSE and the Ensemble Spread of Y_1 and Y_2 are plotted in Fig. 10. It is
14 shown shows that, these two indicators gradually decrease as assimilation step
15 increases and finally reach a low value finally, suggesting that the estimations of Y_1
16 and Y_2 in this case are accurate and effective. Meanwhile, the difference between the
17 RMSE and the Ensemble Spread is small, indicating that the SEOD method estimates
18 the uncertainty properly.

19 [Figure 10]

20 To evaluate the performance of data match and model prediction, the initial and
21 final ensembles of Y_1 and Y_2 are taken into the synthetic model respectively to
22 calculate the head and concentration data, which are then compared with the real

1 observations. It should be noted that only three wells, (black circles in Fig. 7), are
2 selected randomly from the study domain to show the evaluation results of data match
3 and model prediction due to the space limitation. As shown in Fig. 11 and 12, the data
4 calculated by the final ensemble are very close to the real observations, performing
5 much better than those calculated by the initial ensemble. It indicates that the
6 estimated fields of Y_1 and Y_2 are both of low uncertainty after assimilating the most
7 informative measurements from the optimal sampling locations.

8 [Figure 11]

9 [Figure 12]

10 **4. Discussion**

11 **4.1 Comparison of sequential optimization strategy and conventional strategy**

12 In order to illustrate and demonstrate the effectiveness and efficiency of the
13 SEOD method for jointly estimating physical and geochemical parameters, the
14 sequential optimization strategy is compared with the conventional strategy in this
15 subsection. For the convenience of comparison, several synthetic cases (Case 11, 12
16 13) are constructed based on Case 2 by replacing the sequential optimization strategy
17 with the conventional strategy (different fixed sampling locations numbers for
18 different cases). Except this, the other model parameters of cases constructed here are
19 the same as those of Case 2. More details are given in Table 2.

20 Fig.13 shows the *RMSE* and the *Ensemble Spread* for different cases. It
21 illustrates that the sequential optimization strategy obtains better performance of the
22 parameter estimation when the number of sampling locations is the same. Even more,

1 the sequential optimization strategy with 2 optimal sampling locations (Case 2)
2 performs better than the conventional strategy with 10 fixed sampling locations (Case
3 12). Besides, the conventional strategy with a large number of fixed sampling
4 locations could result in the *Ensemble Spread* becoming very small at the first few
5 assimilation steps, which could prevent assimilating further measurements.

6 [Figure 13]

7 4.1-2 Effect of ensemble size

8 All results shown so far until now of Case 2 are based on an ensemble of 100
9 realizations. To evaluate the impact of the ensemble size on the parameter estimation,
10 an analysis with an ensemble of 50, 300, 500, 1000 realizations (Table 2) is performed
11 here in this subsection.

12 The RMSE and the *Ensemble Spread* of different cases ensemble sizes are shown
13 in Fig. 13-14 below. It is shown shows that an appropriate ensemble size is important
14 for the parameter estimation. If the ensemble size is too small (Case 3), ensemble
15 collapse, a phenomenon—that in which the *Ensemble Spread* is artificially small
16 relative to its RMSE, could happen. If the ensemble size is too large (Case 5, 6), it
17 could lead to more computational burden and introduce more observation errors into
18 the model as the SEOD method is based on the Monte Carlo method. It is shown
19 shows that the RMSE and the *Ensemble Spread* of Y_1 and Y_2 are small and close to
20 each other when the ensemble size is 100, suggesting that the estimations of Y_1 and Y_2
21 are accurate and the model uncertainty is estimated properly. Accordingly, the
22 ensemble size is set to 100 in the following cases discussed below.

1 [Figure 43|4]

2 **4.2-3 Effect of the number of optimal sampling locations**

3 Optimizing too many sampling locations could bring a heavy computational
4 burden. Here, to explore the impact of the number of optimal sampling locations on
5 the parameter estimation is explored, several synthetic cases with different numbers
6 of optimal sampling locations are constructed. More details are given in Table 2.

7 As shown in Fig. 44-15 below, the RMSE is no longer sensitive to the number
8 of optimal sampling locations any more when the number of optimal sampling
9 locations is large enough, suggesting that there could be a threshold value of the
10 number of optimal sampling locations in this synthetic model. On the one hand
11 Besides, if the number of optimal sampling locations is too large, the Ensemble
12 Spread becomes extremely small at the first few assimilation steps, which could
13 prevent the method from assimilating further measurements into the model. On the
14 other hand In addition, too many optimal sampling locations could lead to high
15 economic cost and heavy computational burden. Therefore, 2 to 5 optimal sampling
16 locations are enough and appropriate in this model.

17 [Figure 44|5]

18 **4.3 Comparison between sequential optimization strategy and conventional**
19 **strategy**

20 In order to demonstrate the effectiveness and efficiency of the SEOD method for
21 jointly estimating physical and geochemical parameters, sequential optimization
22 strategy is compared with conventional strategy in this subsection. More details are

1 given in Table 2.

2 The RMSE and Ensemble Spread of different cases are plotted in Fig. 15 below.

3 It is intuitively seen that sequential optimization strategy obtains better performance

4 of parameter estimation when the number of sampling locations is the same. Even

5 more, sequential optimization strategy with 2 optimal sampling locations (Case 2)

6 performs better than conventional strategy with 10 fixed sampling locations (Case 12).

7 Besides, conventional strategy with large number of fixed sampling locations could

8 lead that the Ensemble Spread becomes very small at first few assimilation steps,

9 which could prevent assimilating further measurements.

10 [Figure 15]

11 4.4 Improvement of the SEOD method

12 In the original SEOD method, since the EnKF is a sequential history matching

13 method, the optimization algorithm part (GA) needs to be invoked N_s times to obtain

14 the optimal sampling design at each assimilation step, which is time-consuming. To

15 enhance its computational efficiency, we improve the original SEOD method by

16 replacing the EnKF with the ES-MDA. The loop of the improved SEOD method is

17 shown in Fig. 16.

18 [Figure 16]

19 Fig. 16 shows that the loop of the improved SEOD method is divided into outer

20 loop and inner loop. The outer loop is similar to the original SEOD method, which

21 consists of a forecast step, an optimal design step and an analysis step. The inner loop

22 of the improved SEOD method is part of the ES-MDA. Unlike the EnKF, the

1 ES-MDA performs N_a times global update so as to assimilate the same data (all
 2 available data) multiple times without restarting the forward model, which helps
 3 enhance the computational efficiency. In the improved SEOD method, we divide all
 4 N_s assimilation steps in the original SEOD method into N_g groups, the following loop
 5 is performed for each group in chronological order (from 1 to N_g). Note that N_a of all
 6 cases in this subsection (Case 14, 15, 16, 17) is set to 4 (Emerick and Reynolds 2013).
 7 More details are given in Table 2.

8 (I) Forecast step

9 Run the forward model G from beginning time step of the group $j+1$ to the end
 10 time step of the group $j+1$ with updated parameters from the group j (Eq. (10)).

11
$$\mathbf{x}_{i,j+1}^f = G(\mathbf{x}_{i,j}^a), \quad i=1,2,\dots,N_e \quad (10)$$

12 In the above equation, i is the ensemble member index, j is the group index (from
 13 1 to N_g), superscripts f and a denote forecast and analysis, respectively.

14 (II) Optimal design

15 This step is similar with the original SEOD method, using the GA to solve an
 16 optimization problem to obtain the most informative measurements from the optimal
 17 sampling design.

18 (III) Analysis step

19 The following update equation (Eq. (11)) of the ES-MDA is different from that
 20 of the EnKF.

21
$$\mathbf{x}_{i,j+1}^a = \mathbf{x}_{i,j+1}^f + \mathbf{C}_{YD} (\mathbf{C}_{DD} + \alpha_l \mathbf{C}_D)^{-1} (\mathbf{d}_{obs_{i,j+1}} - \mathbf{d}_{i,j+1}), \quad i=1,2,\dots,N_e \quad (11)$$

22 In the above equation, l is the times index of the ES-MDA, $l=1,2,\dots,N_a$; \mathbf{d}_{obs}

1 is the perturbed observations with noise of covariance $\alpha_l \mathbf{C}_D$ ($\alpha_1=9.333$, $\alpha_2=7.0$,
2 $\alpha_3=4.0$ and $\alpha_4=2.0$, Emerick and Reynolds 2013). Other letters in this equation have
3 the same meaning as those in Eq. (6).

4 After N_a times global update, the updated ensemble X^a of the group $j+1$ is
5 obtained here. Then, go back to step (I), the updated ensemble is implemented for the
6 next group. Through this improvement, the times of invoking the GA decrease from
7 N_s to N_g , which helps enhance the computational efficiency.

8 To compare the improved SEOD method with the original one and discuss the
9 influence factors of the improved one, several cases are constructed with all model
10 parameters the same as those in Case 2. More details can be found in Table 2. The
11 results of these cases are shown in Fig. 17.

12 [Figure 17]

13 Fig. 17 shows that the number of optimal sampling locations dominantly affects
14 the results of the improved SEOD method. In Case 2, two optimal sampling locations
15 are chosen at each assimilation step. During the whole data assimilation, total 20
16 sampling locations are used to obtain measurements at most (2×10 , if the optimal
17 sampling locations are different at each step). Therefore, if the number of optimal
18 sampling locations in the improved SEOD method is too small, the improved SEOD
19 method can't estimate the parameters of the entire study domain well just through a
20 few sampling locations (e.g., Case 14 with total 4 optimal sampling locations at most).
21 When the number of sampling locations is enough, the result of the improved SEOD
22 method is acceptable (e.g., Case 15 with total 10 optimal sampling locations at most).

1 The result of Case 15 is comparable with the result of Case 2, but the computer cost of
2 Case 15 is much less than that of Case 2 because that Case 15 just invokes the GA
3 only twice while Case 2 invokes the GA 10 times. Therefore, if the number of optimal
4 sampling locations is not set to a very small value, the computational efficiency will
5 be enhanced by using the improved SEOD method.

6 Furthermore, the way of dividing the observation time (assimilation steps in the
7 original SEOD method) into several groups (called the group division strategy in the
8 following context) affects the results of the improved method as well. From the
9 comparison of Case 15, Case 16 and Case 17, it is obvious that Case 16, whose
10 observation time in each divided groups is progressively increasing, has a better data
11 assimilation result than the other two cases, whose observation time in each divided
12 groups is equivalent or progressively decreasing. It is an interesting phenomenon,
13 which is worth further research. For now, we think it is probably because more and
14 more precise measurements in the early stage of the data assimilation would lead to
15 excessive update of estimated parameters (Burgers et al. 1998; Evensen 2009).
16 Therefore, future studies should focus on optimizing the group division strategy to
17 obtain a more accurate estimation of model parameters.

18 **5. Conclusions**

19 In this study, the effectiveness and efficiency of the sequential ensemble based
20 optimal design (SEOD) method for jointly estimating the spatial distribution of
21 physical and geochemical parameters are illustrated and demonstrated by two
22 synthetic cases. The results indicate that uncertainties of both physical and

1 geochemical parameters decrease after assimilating the most informative
2 measurements from optimal sampling locations.

3 Ensemble size and number of optimal sampling locations have impact on the
4 parameter estimation. For ensemble size, too small ensemble size could lead to
5 ensemble collapse, and heavier computational burden and more observation errors
6 could be caused conversely. For the number of optimal sampling locations, when it is
7 large enough, the *RMSE* is not sensitive to the number of optimal sampling locations
8 any more. Overall, appropriate ensemble size and number of optimal sampling
9 locations can estimate parameters well with low computational burden and economic
10 cost.

11 Compared with conventional strategy, sequential optimization strategy has better
12 performance of parameter estimation when the number of sampling locations is the
13 same. Even more, sequential optimization strategy with 2 optimal sampling locations
14 (Case 2) performs better than conventional strategy with 10 fixed sampling locations
15 (Case 12).

16 In addition, only two kinds of measurements, head and concentration, are
17 assimilated in this work. For further study, more kinds of measurements, such as
18 temperature, geophysical data and so on, can be taken into data assimilation so as to
19 make the most of both solid and soft measurements to improve the accuracy of
20 parameter estimation. In this study, we make use of a sequential ensemble-based
21 optimal design (SEOD) method to jointly estimate physical and geochemical
22 parameters of groundwater models.

1 Both physical and geochemical parameters are estimated accurately in the
2 one-dimensional and two-dimensional synthetic cases by using the SEOD method.
3 Uncertainties of both physical and geochemical parameters decrease after assimilating
4 the most informative measurements at the optimal sampling locations, and the
5 accuracy of model prediction increase meanwhile. Furthermore, several comparison
6 cases are tested and analyzed, results illustrate and demonstrate the effectiveness and
7 efficiency of the SEOD method on jointly estimating high-dimensional physical and
8 geochemical parameters in groundwater models.

9 The ensemble size and the number of optimal sampling locations have impacts
10 on the parameter estimation based on the SEOD method. A too small ensemble size
11 would lead to the ensemble collapse. Furthermore, when the number of optimal
12 sampling locations is too large, heavier computational burden and more observation
13 errors would be caused, and the *RMSE* is no longer sensitive to the number of optimal
14 sampling locations. How to determine the optimal ensemble size and sampling
15 locations number for different scenarios is worth further investigation.

16 The original SEOD method has a heavy computational burden because it invokes
17 the GA too many times. To enhance its computational efficiency, we proposed an
18 improved SEOD method in this study by replacing the EnKF with the ES-MDA. The
19 results of comparison cases show that the improved SEOD method is advantageous
20 than the original one, which makes the SEOD framework more promising for the
21 parameter estimation and the optimal sampling strategy design. The number of
22 optimal sampling locations and the strategy of dividing groups would affect the

1 results of the improved SEOD method.

2 It is noted that only two kinds of measurements (head and concentration) are
3 assimilated in this work. More kinds of measurements (e.g., hydraulic conductivity,
4 porosity, temperatures and hydrogeophysical data) can be assimilated simultaneously
5 so as to make use of more hard and soft data to improve the accuracy of parameter
6 estimation in further study.

7

8 Acknowledgements

9 The authors would like to thank the anonymous referees for their insightful comments and
10 suggestions that have helped improve the paper. This work was financially supported by
11 the National Nature Science Foundation of China grants (No. U1503282, 41672229,
12 and 41172206). We would like to thank Mr. Jun Man from Zhejiang University for
13 providing the SEOD code.

14

15 References

- 16 Aanonsen SI, Nævdal G, Oliver DS, Reynolds AC, Vallès B (2009) The ensemble Kalman filter in
17 reservoir engineering--a review. *Spe Journal* 14:393–412
- 18 Atchley AL, Navarre-Sitchler AK, Maxwell RM (2014) The effects of physical and geochemical
19 heterogeneities on hydro-geochemical transport and effective reaction rates. *J Contam
Hydrol* 165:53–64 doi:10.1016/j.jconhyd.2014.07.008
- 20 Bear J (1972) Dynamics of Fluids in Porous Materials. Dover, New York
- 21 Burgers G, Leeuwen P, Evensen G (1998) Analysis scheme in the ensemble Kalman filter. *Month
Weather Rev* 126:1719–1724
- 22 Carrera J, Alcolea A, Medina A, Hidalgo J, Slooten LJ (2005) Inverse problem in hydrogeology.
23 *Hydrogeology Journal* 13(1):206–222
- 24 Carrera J, Neuman SP (1986) Estimation of Aquifer Parameters Under Transient and Steady State
25 Conditions: 3. Application to Synthetic and Field Data. *Water Resources Research*
26 22:228–242
- 27 Catania F, Massabó M, Paladino O (2004) Optimal Sampling Strategy for Parameters Estimation.

- 1 **In: AGU Fall Meeting, 2004**
- 2 Chen Y, Oliver DS (2010) Cross-covariances and localization for EnKF in multiphase flow data
3 assimilation. Computational Geosciences 14:579-601
- 4 Chen Y, Zhang D (2006) Data assimilation for transient flow in geologic formations via ensemble
5 Kalman filter. Advances in Water Resources 29:1107-1122
6 doi:10.1016/j.advwatres.2005.09.007
- 7 **Cleveland TG, Yeh WG (1991) Optimal Configuration and Scheduling of Ground Water Tracer**
8 **Test. Journal of Water Resources Planning & Management 117:37-51**
- 9 Cleveland TG, Yeh WWG (1990) Sampling network design for transport parameter identification.
10 Journal of Water Resources Planning & Management 116:764-783
- 11 Dagan G (1984) Solute transport in heterogeneous porous formations. Journal of Fluid Mechanics
12 145:151 doi:10.1017/s0022112084002858
- 13 Dagan G (1985) Stochastic modeling of groundwater flow by unconditional and conditional
14 probabilities: The inverse problem. Water Resources Research 21(1):65-72
- 15 Doherty J (2004) PEST: Model-Independent Parameter Estimation,User's Manual (5th edition).
16 Watermark Numerical Computing, Australia
- 17 Emerick AA, Reynolds AC (2011) Combining sensitivities and prior information for covariance
18 localization in the ensemble Kalman filter for petroleum reservoir applications.
19 Computational Geosciences 15:251-269
- 20 Emerick AA, Reynolds AC (2013) Ensemble smoother with multiple data assimilation. Comput
21 Geosci-Uk 55:3-15 doi:10.1016/j.cageo.2012.03.011
- 22 Evensen G (2003) The Ensemble Kalman Filter: theoretical formulation and practical
23 implementation. Ocean Dynamics 53:343-367
- 24 Evensen G (2009) Data assimilation: the ensemble Kalman filter. Springer, Berlin
- 25 Fennell DE, Carroll AB, Gossett JM, Zinder SH (2001) Assessment of indigenous reductive
26 dechlorinating potential at a TCE-contaminated site using microcosms, polymerase chain
27 reaction analysis, and site data. Environmental Science & Technology 35:1830-1839
- 28 Gómez-Hernández JJ, Hendricks Franssen HJ, Sahuquillo A (2003) Stochastic conditional inverse
29 modeling of subsurface mass transport: A brief review and the self-calibrating method.
30 Stochastic Environmental Research and Risk Assessment 17(5):319-328
- 31 Gu Y, Oliver DS (2007) An Iterative Ensemble Kalman Filter for Multiphase Fluid Flow Data
32 Assimilation. Spe Journal 12:1990 - 1995
- 33 Harbaugh AW, Banta ER, Hill MC, McDonald MG (2000) MODFLOW-2000, The U. S.
34 Geological Survey Modular Ground-Water Model—User Guide to Modularization
35 Concepts and the Ground-Water Flow Process. U.S. Geological Survey Open-File Report
36 00-92, 121 p
- 37 Hendricks Franssen HJ, Alcolea A, Riva M, Bakr M, van der Wiel N, Stauffer F, Guadagnini A
38 (2009) A comparison of seven methods for the inverse modelling of groundwater flow.
39 Application to the characterisation of well catchments. Advances in Water Resources
40 32(6):851-872
- 41 Hendricks Franssen HJ, Kinzelbach W (2008) Real-time groundwater flow modeling with the
42 Ensemble Kalman Filter: Joint estimation of states and parameters and the filter
43 inbreeding problem. Water Resour Res 44:354-358
- 44 Huang C, Hu BX, Li X, Ye M (2009) Using data assimilation method to calibrate a heterogeneous

- 1 conductivity field and improve solute transport prediction with an unknown
2 contamination source. Stochastic Environmental Research and Risk Assessment
3 23(8):1155
- 4 Kalman RE (1960) A new approach to linear filtering and prediction problems. Trans ASME J
5 Basic Eng 82(D):35–45
- 6 Knopman DS, Voss CI (1987) Behavior of sensitivities in the one-dimensional
7 advection-dispersion equation: Implications for parameter estimation and sampling design.
8 Water Resources Research 23:253–272
- 9 ~~Knopman DS, Voss CI (1988) Further comments on sensitivities, parameter estimation, and
10 sampling design in one dimensional analysis of solute transport in porous media. Water
11 Resources Research 24:225–238~~
- 12 ~~Knopman DS, Voss CI (1989) Multiobjective sampling design for parameter estimation and model
13 discrimination in groundwater solute transport. Water Resources Research 25:2245–2258~~
- 14 ~~Knopman DS, Voss CI, Garabedian SP (1991) Sampling design for groundwater solute transport:
15 Tests of methods and analysis of Cape Cod tracer test data. International Journal of Rock
16 Mechanics & Mining Sciences & Geomechanics Abstracts 28:925–949~~
- 17 Kullback S (1997) Information theory and statistics. Courier Corporation
- 18 Li G, Reynolds AC (2009) Iterative Ensemble Kalman Filters for Data Assimilation. Spe Journal
19 14:496–505
- 20 Li L, Steefel CI, Kowalsky MB, Englert A, Hubbard SS (2010) Effects of physical and
21 geochemical heterogeneities on mineral transformation and biomass accumulation during
22 biostimulation experiments at Rifle, Colorado. J Contam Hydrol 112:45–63
23 doi:10.1016/j.jconhyd.2009.10.006
- 24 Man J, Zhang J, Li W, Zeng L, Wu L (2016) Sequential ensemble-based optimal design for
25 parameter estimation. Water Resour Res 52:7577–7592 doi:10.1002/2016wr018736
- 26 Moradkhani H, Sorooshian S, Gupta HV, Houser PR (2005) Dual state-parameter estimation of
27 hydrological models using ensemble Kalman filter. Advances in Water Resources
28 28:135–147
- 29 Neuman SP (1973) Calibration of distributed parameter groundwater flow models viewed as a
30 multiple objective decision process under uncertainty. Water Resources Research
31 9(4):1006–1021
- 32 ~~Nishikawa T, Yeh WG (1989) Optimal pumping test design for the parameter identification of
33 groundwater systems. Water Resources Research 25:1737–1747~~
- 34 Nowak W, De Barros FPJ, Rubin Y (2010) Bayesian geostatistical design: Task-driven optimal site
35 investigation when the geostatistical model is uncertain. Water Resources Research 46(3):
36 374–381 doi:10.1029/2009WR008312
- 37 Oliver DS, Chen Y (2011) Recent progress on reservoir history matching: a review. Computational
38 Geosciences 15(1):185–221
- 39 Oliver DS, Cunha LB, Reynolds AC (1997) Markov chain Monte Carlo methods for conditioning
40 a permeability field to pressure data. Mathematical Geology 29(1): 61–91
- 41 Prommer CH, Post V (2010) A Reactive Multicomponent Transport Model for Saturated Porous
42 Media. Groundwater 48(5):627–632
- 43 Rubin Y (1991) Transport in heterogeneous porous media: Prediction and uncertainty. Water
44 Resour Res 27:1723–1738

- 1 Sandrin SK, Brusseau ML, Piatt JJ, Bodour AA, Blanford WJ, Nelson NT (2004) Spatial
2 variability of in situ microbial activity: biotracer tests. *Groundwater* 42:374-383
- 3 Scheibe TD et al. (2006) Transport and biogeochemical reaction of metals in a physically and
4 chemically heterogeneous aquifer. *Geosphere* 2:220-235 doi:10.1130/Ges00029.1
- 5 Scheibe TD, Fang Y, Murray CJ, Roden EE, Chen J, Chien YJ, Brooks SC, Hubbard SS (2006)
6 Transport and biogeochemical reaction of metals in a physically and chemically
7 heterogeneous aquifer. *Geosphere* 2(4):220-235 doi:10.1130/Ges00029.1
- 8 Sorensen JVT, Madsen H, Madsen H (2004) Data assimilation in hydrodynamic modelling: on the
9 treatment of non-linearity and bias. *Stoch Environ Res Risk Assess* 18(7):228-244
- 10 Sudicky EA (1986) A natural gradient experiment on solute transport in a sand aquifer: Spatial
11 variability of hydraulic conductivity and its role in the dispersion process. *Water Resour
12 Res* 22:2069-2082 doi:10.1029/WR022i013p02069
- 13 Sun NZ, Yeh WWG (2007) Development of objective-oriented groundwater models: 2. Robust
14 experimental design. *Water resources research* 43(2) doi:10.1029/2006WR004888
- 15 Tong J, Hu BX, Yang J (2010) Using data assimilation method to calibrate a heterogeneous
16 conductivity field conditioning on transient flow test data. *Stochastic environmental
17 research and risk assessment* 24(8):1211-23
- 18 Ushijima TT, Yeh WWG (2015) Experimental design for estimating unknown hydraulic
19 conductivity in an aquifer using a genetic algorithm and reduced order model. *Advances
20 in Water Resources* 86:193-208
- 21 Van Leeuwen PJ, Evensen G (1996) Data Assimilation and Inverse Methods in Terms of a
22 Probabilistic Formulation. *Monthly Weather Review* 124:2898-2913
- 23 Whitley D (1994) A genetic algorithm tutorial. *Stat Comput* 4(2):65-85
- 24 Zhang D, Lu Z (2004) An efficient, high-order perturbation approach for flow in random porous
25 media via Karhunen-Loève and polynomial expansions. *Journal of Computational
26 Physics* 194:773-794
- 27 Zhang J, Zeng L, Chen C, Chen D, Wu L (2015) Efficient Bayesian experimental design for
28 contaminant source identification. *Water Resources Research* 51(1):576-598
- 29 Zheng C (2006) MT3DMS v5.2 supplemental user's guide: Technical report to the U.S. Army
30 Engineer Research and Development Center, Department of Geological Sciences,
31 University of Alabama, p 24
- 32 Zhou HY, Gómez-Hernández JJ, Franssen HJJ, Li LP (2011) An approach to handling
33 non-Gaussianity of parameters and state variables in ensemble Kalman filtering. *Adv
34 Water Resour* 34:844-864 doi:10.1016/j.advwatres.2011.04.014
- 35 Zhou HY, Gómez-Hernández JJ, Li LP (2014) Inverse methods in hydrogeology: Evolution and
36 recent trends. *Advances in Water Resources* 63:22-37
- 37

Tables

Table 1: Flow and transport parameters used in Case 1

Flow simulation	Transient state
Total simulation time (days)	10
Stress period	1
Time steps	100
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	150
Model width (m)	5
Model height (m)	5
Starting head (m)	100
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
Injection rate per well (m^3/d)	50
Pumping rate per well (m^3/d)	45

1 **Table 2: Data assimilation related parameters used in different cases**

Case name	Dimension	Number of ensemble (N_e)	Number of assimilation step (N_s)	Optimize or not	Number of optimal sampling locations
Case 1	1	300	10	Y	2
Case 2	2	100	10	Y	2
Case 3	2	50	10	Y	2
Case 4	2	300	10	Y	2
Case 5	2	500	10	Y	2
Case 6	2	1000	10	Y	2
Case 7	2	100	10	Y	1
Case 8	2	100	10	Y	5
Case 9	2	100	10	Y	10
Case 10	2	100	10	Y	20
Case 11	2	100	10	N	(2 fixed)
Case 12	2	100	10	N	(10 fixed)
Case 13	2	100	10	N	(20 fixed)
Case 14	2	100	8(50, 50)*	Y	2
Case 15	2	100	8(50, 50)	Y	5
Case 16	2	100	12(20, 30, 50)	Y	5
Case 17	2	100	12(50, 30, 20)	Y	5

2 * The numbers in the parentheses are the group division of observation time, and the number in front of
3 the parentheses is the number of assimilation steps. For example, 8(50, 50) represents that there are 8
4 steps in the assimilation and the observation time is divided into two groups with each group having an
5 observation time of 50 days.

6

1 **Table 3: Flow and transport parameters used in Case 2**

Flow simulation	Transient state
Total simulation time (days)	100
Stress period	1
Time steps	200
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	105
Model width (m)	65
Model height (m)	5
Starting head (m)	50
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
horizontal transverse dispersivity (m)	1
Injection rate per well (m^3/d)	80
Pumping rate per well (m^3/d)	120
TCE injection concentration per well (mg/L)	50

2

3

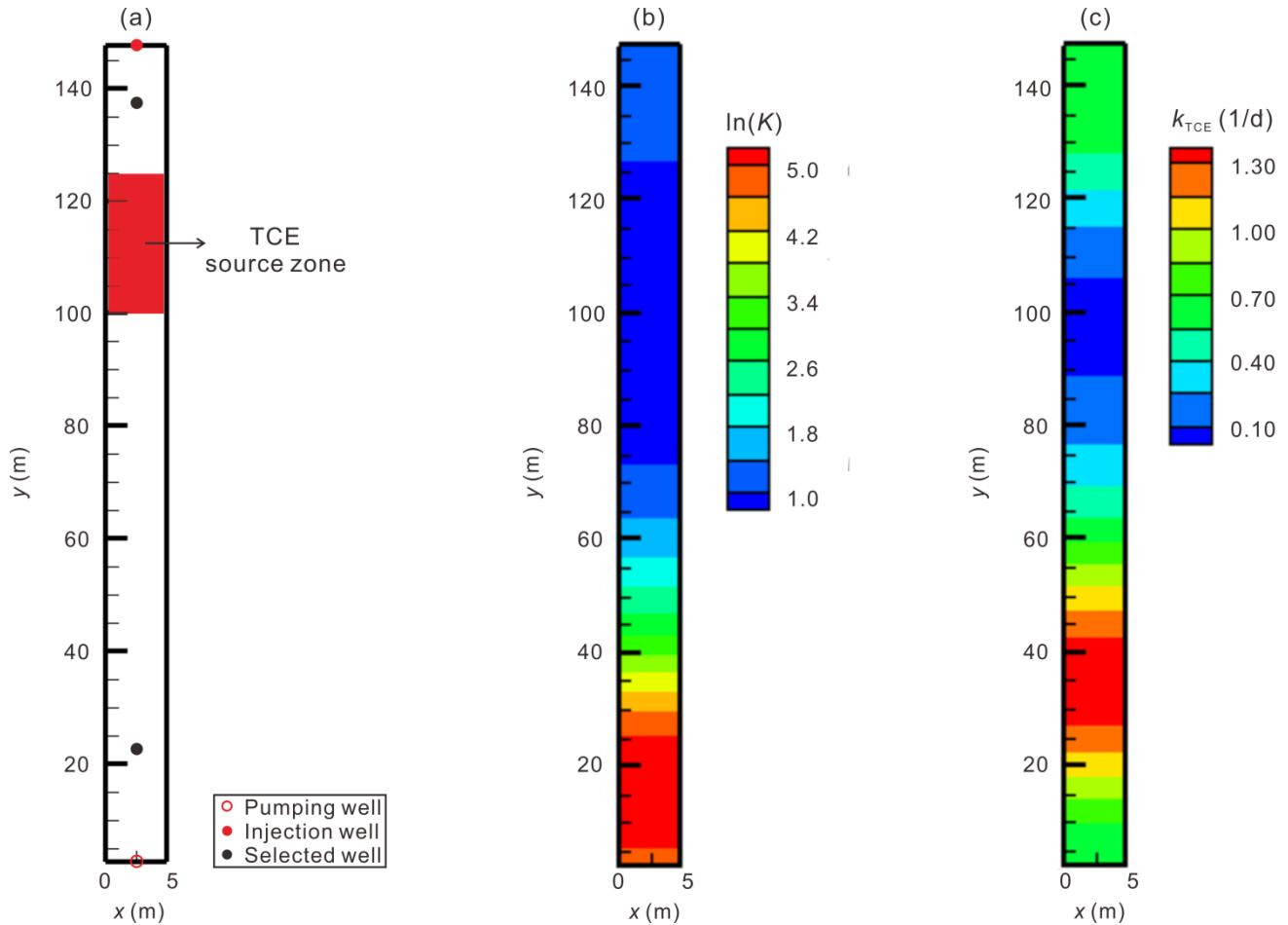
1

Table 4: Computational costs

Case name	Times of model invoking	Times of GA invoking	computing time (using the same computer)
Case 2	10	10	2 h
Case 14	8	2	26 min
Case 15	8	2	26 min
Case 16	12	3	25 min
Case 17	12	3	25 min

2

1 **Figures**



2
3 **Fig. 1** The conceptual model (a), the reference fields of the hydraulic conductivity (b) and the first-order
4 rate constant (c) for Case 1.

5

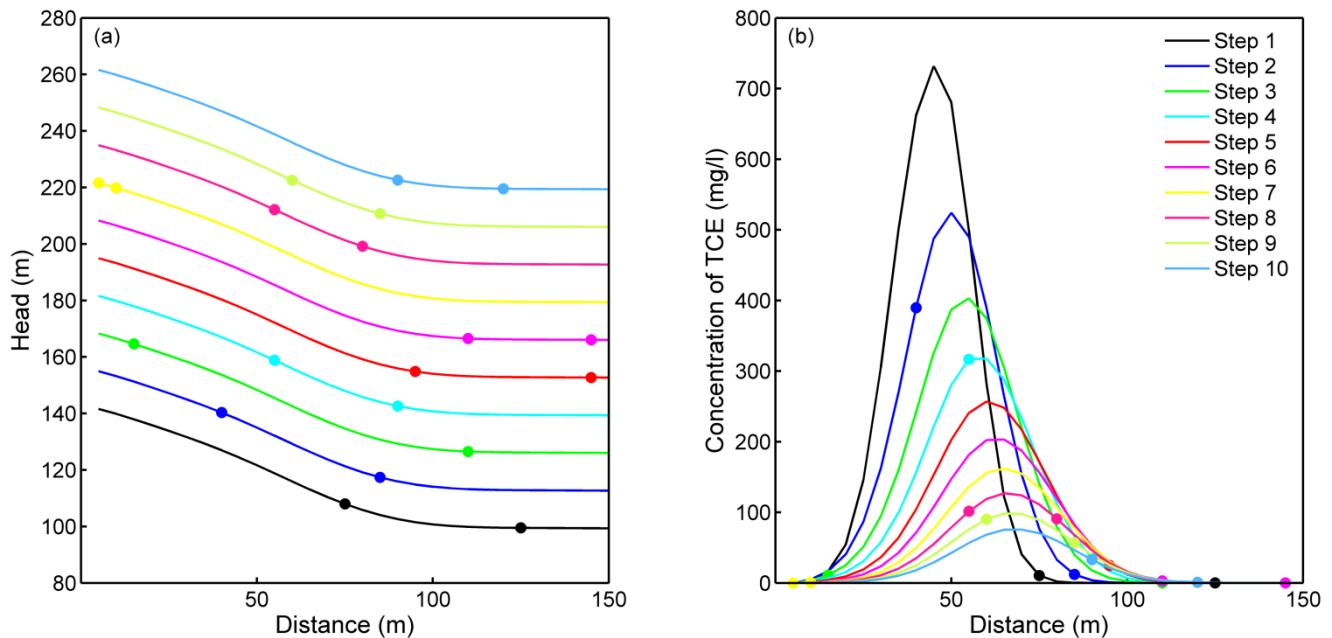
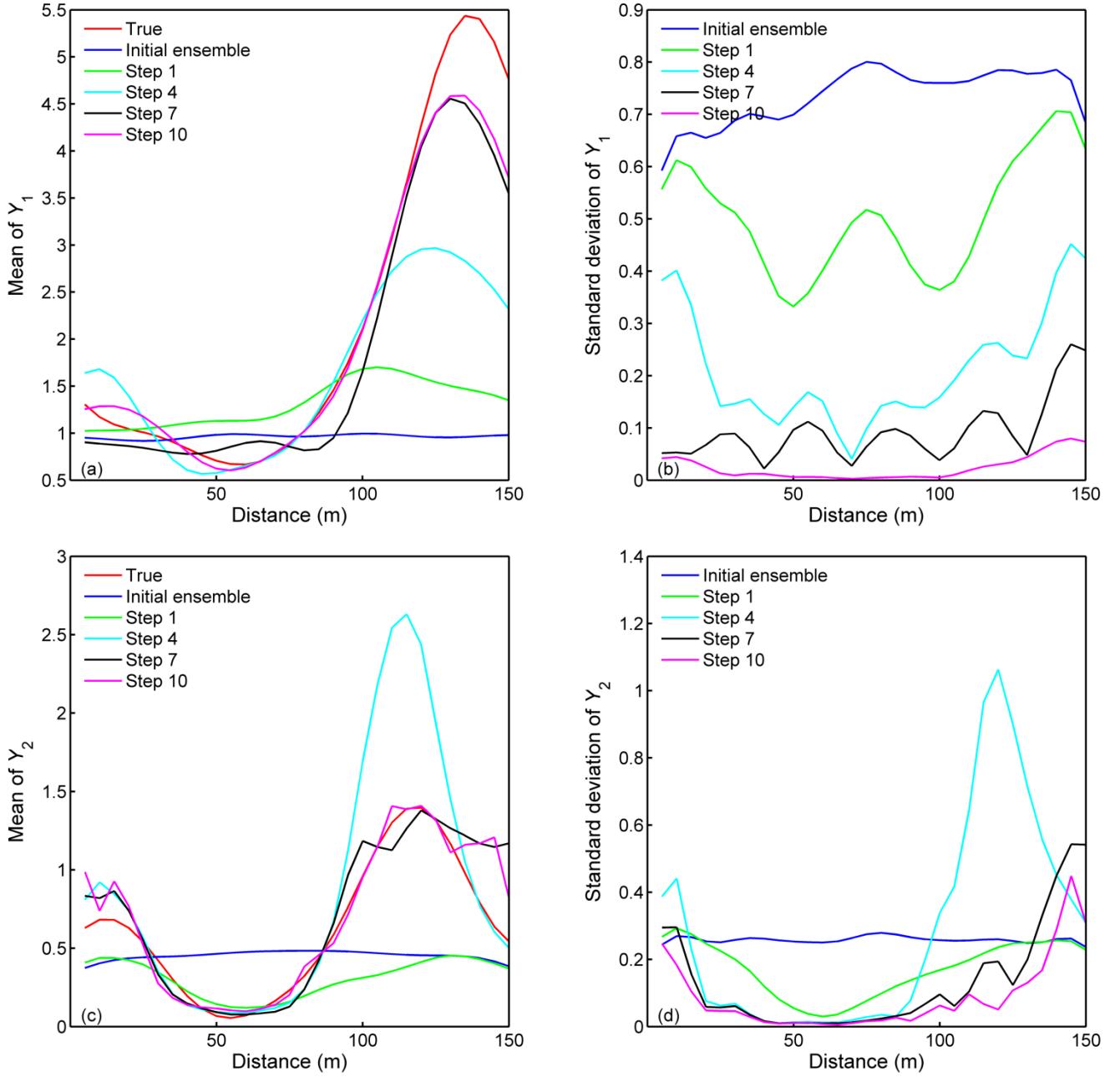


Fig. 2 The calculated flow field (a) and concentration field (b) in Case 1. The circles denote the optimal sampling locations proposed by the SEOD method.



1 **Fig. 3** The ensemble mean and the standard deviation of field Y_1 and Y_2 in Case 1. (a) and (b) are the
2 ensemble mean and the standard deviation of field Y_1 respectively, while (c) and (d) are the ensemble
3 mean and the standard deviation of field Y_2 respectively.
4

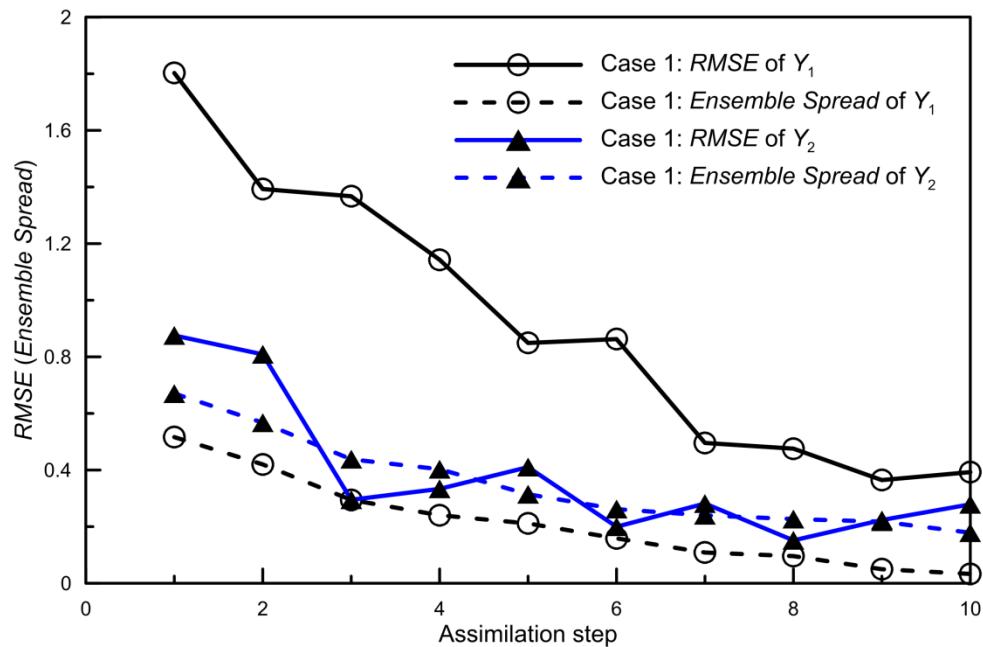


Fig. 4 The performance indicators (the RMSE and the Ensemble Spread) at each assimilation step for field Y_1 and Y_2 in Case 1.

1
2
3
4

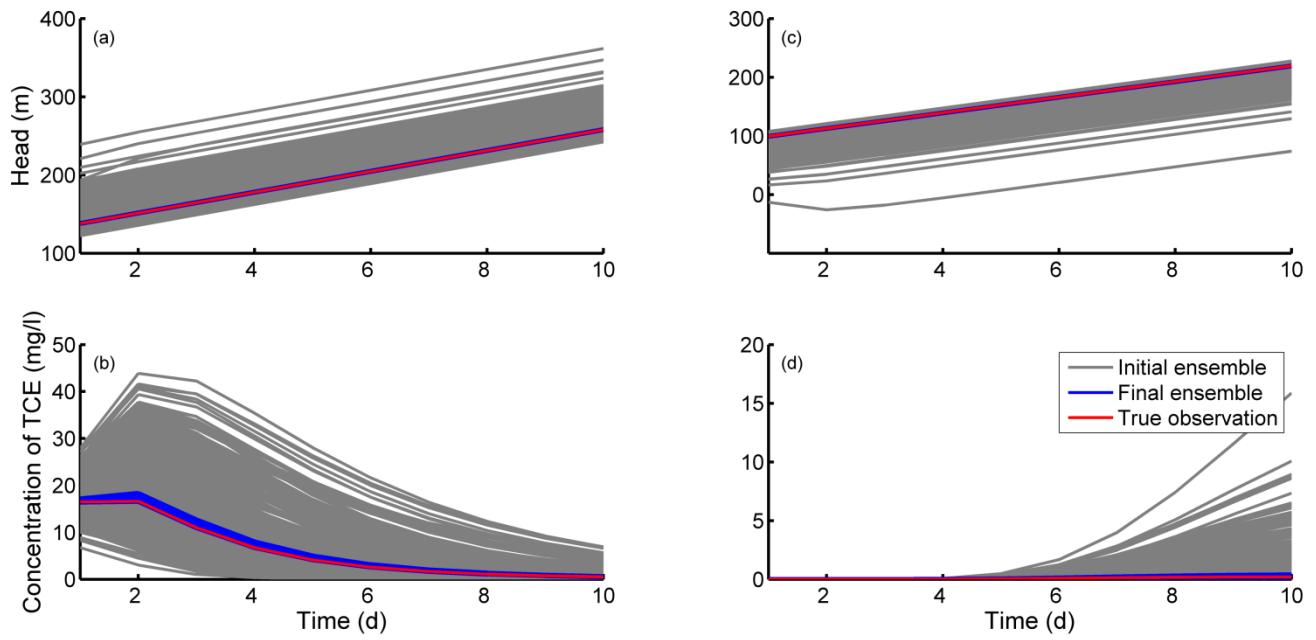


Fig. 5 The performance of data match. (a) and (c) show the data match of the head of two selected wells respectively, while (b) and (d) show the data match of the TCE concentration data of two selected wells respectively. (a) and (b) is the plot of head and TCE concentration data match of well 1 respectively, while (c) and (d) is the plot of head and TCE concentration data match of well 2 respectively.

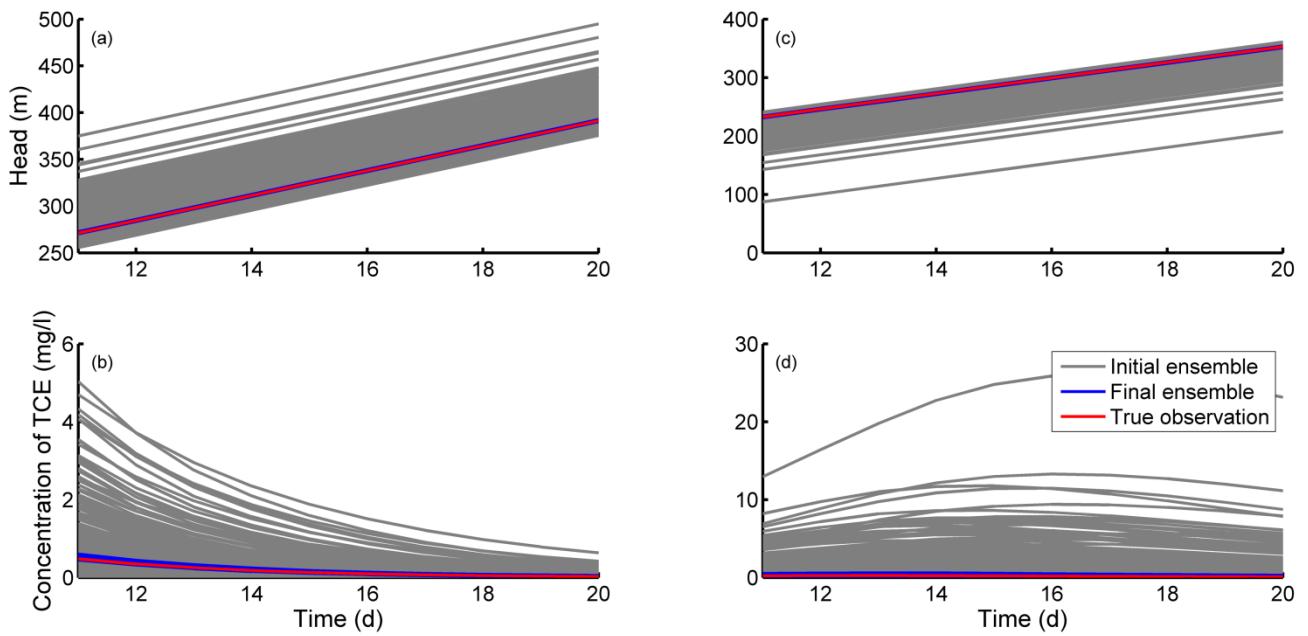


Fig. 6 The performance of model prediction. (a) and (c) show the prediction of the head of two selected wells respectively, while (b) and (d) show the prediction of the TCE concentration data of two selected wells respectively. (a) and (b) is the prediction of head and TCE concentration of well 1 respectively, while (c) and (d) is the prediction of head and TCE concentration data match of well 2 respectively.

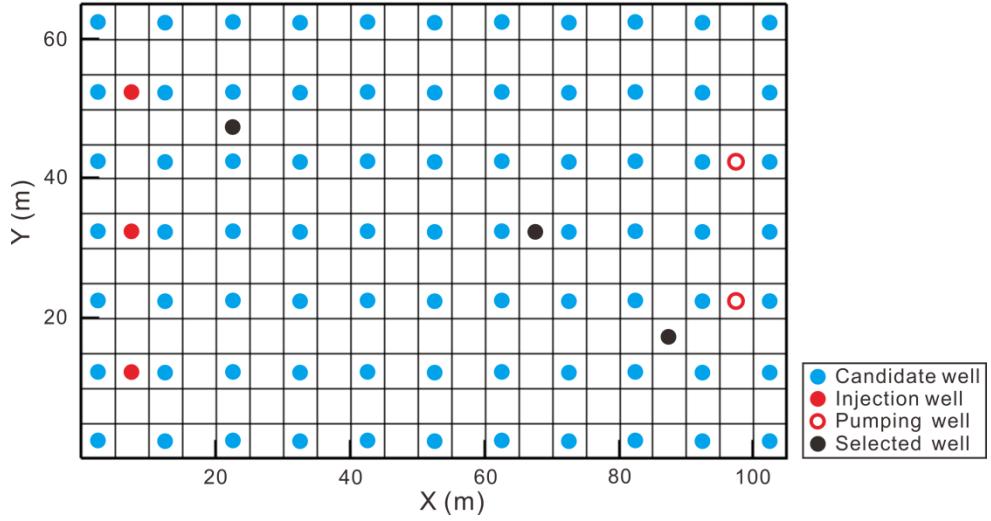
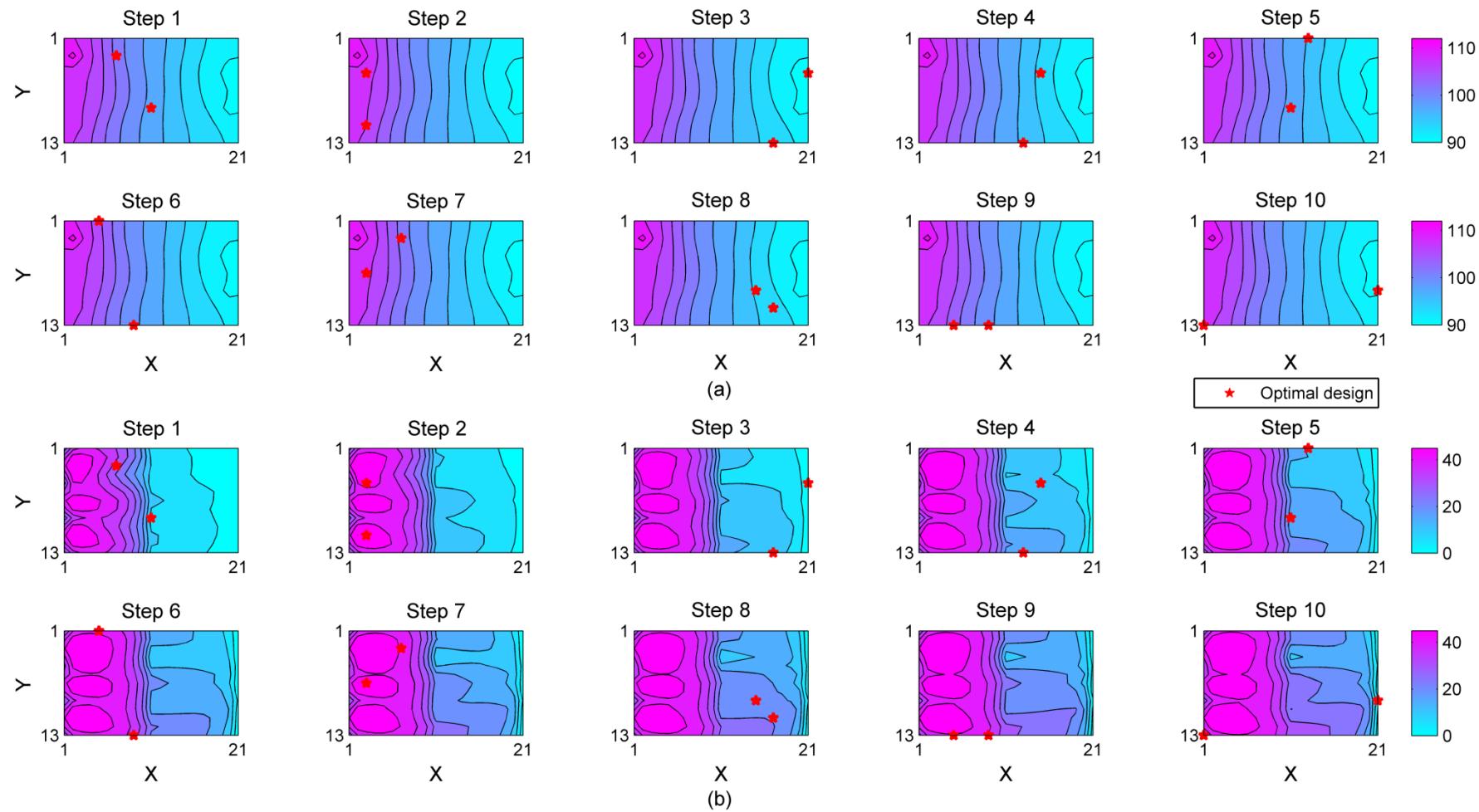


Fig. 7 The schematic of the two-dimensional conceptual model.

1
2
3



1
2 **Fig. 8** The optimal sampling locations (red stars) at every assimilation step of Case 2. (a) and (b) are the contour maps of the flow filed and the concentration
3 field, respectively.
4

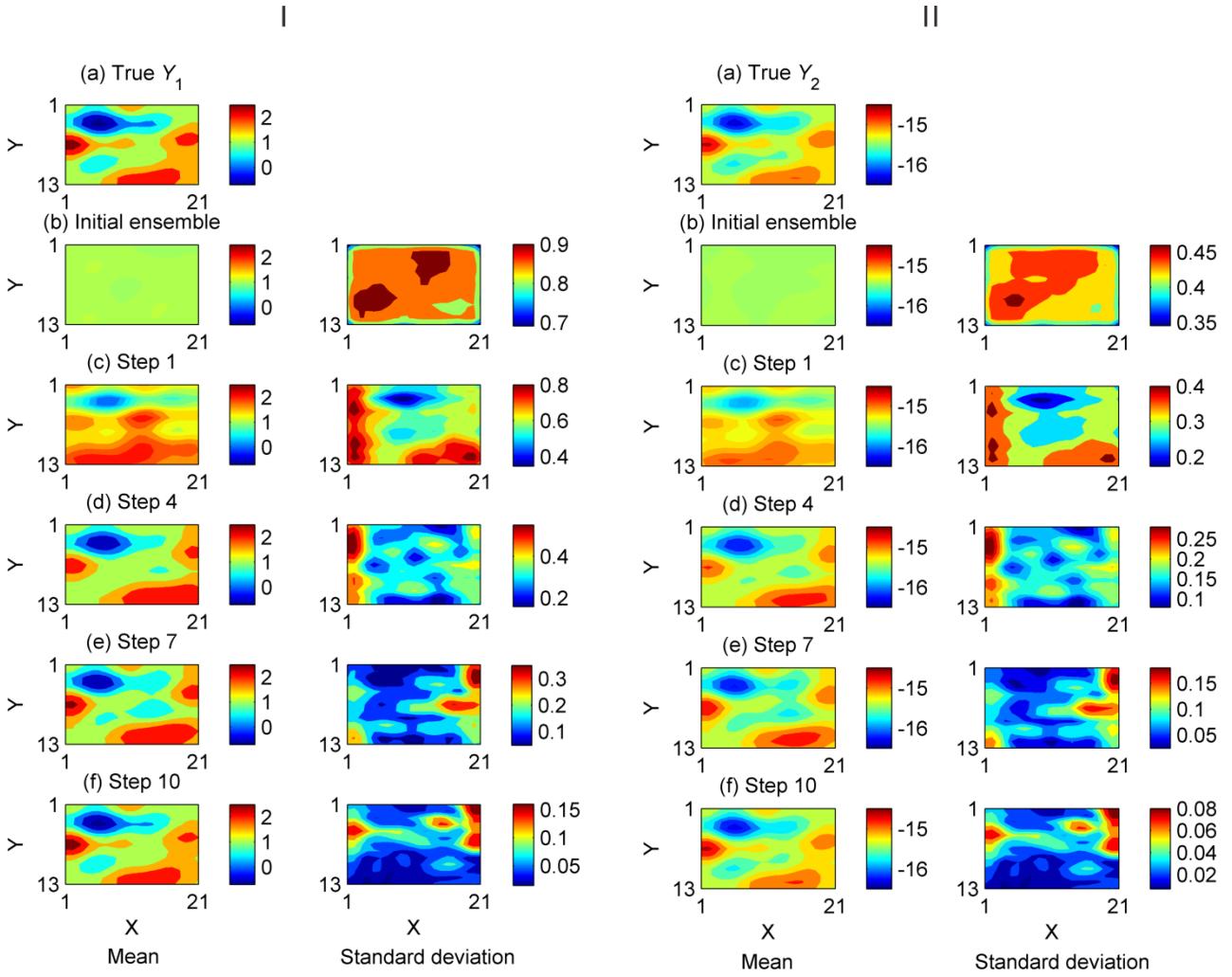


Fig. 9 The ensemble mean and the standard deviation of field Y_1 and Y_2 in Case 2. I for field Y_1 , II for field Y_2 .

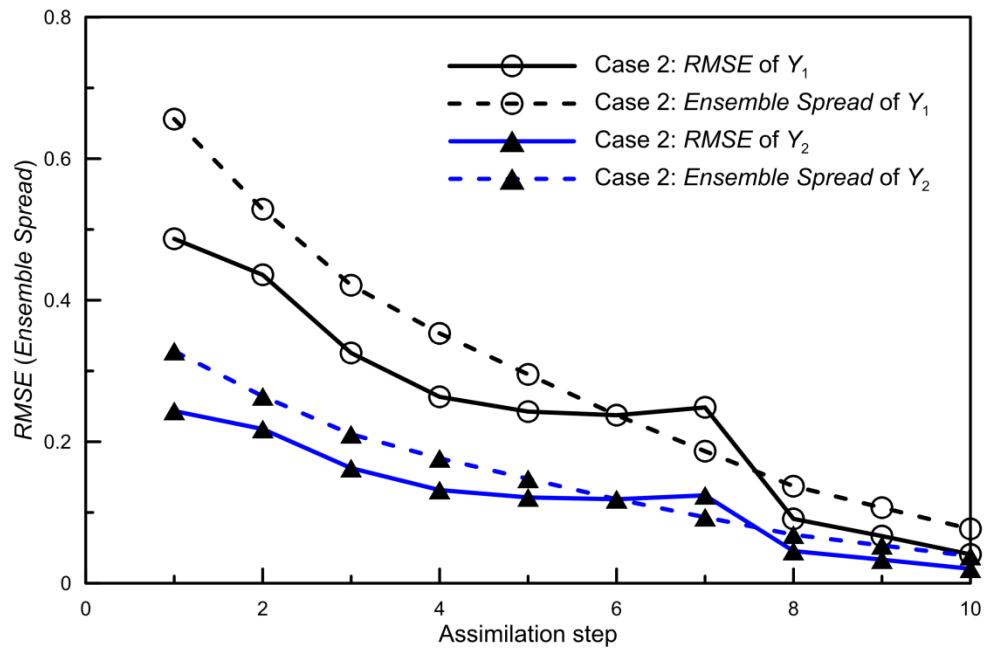


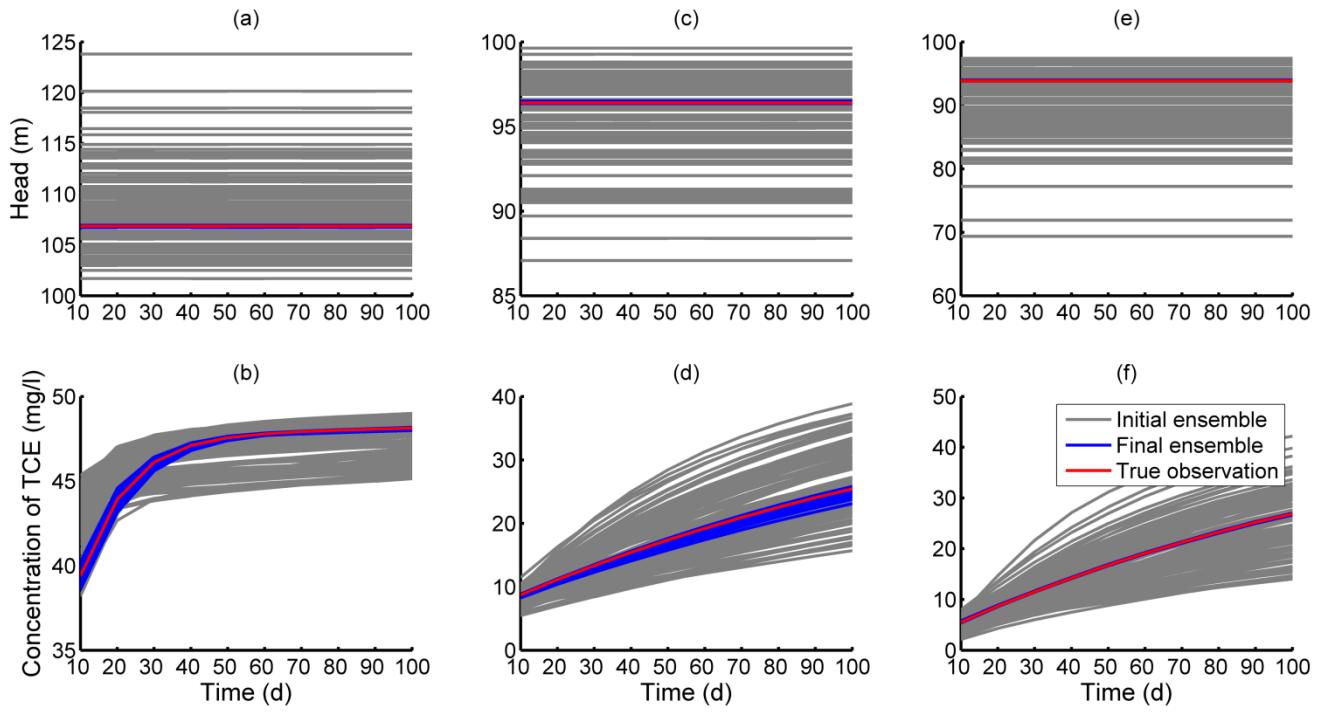
Fig. 10 The performance indicators (the RMSE and the Ensemble Spread) at each assimilation step for field Y_1 and Y_2 in Case 2.

1

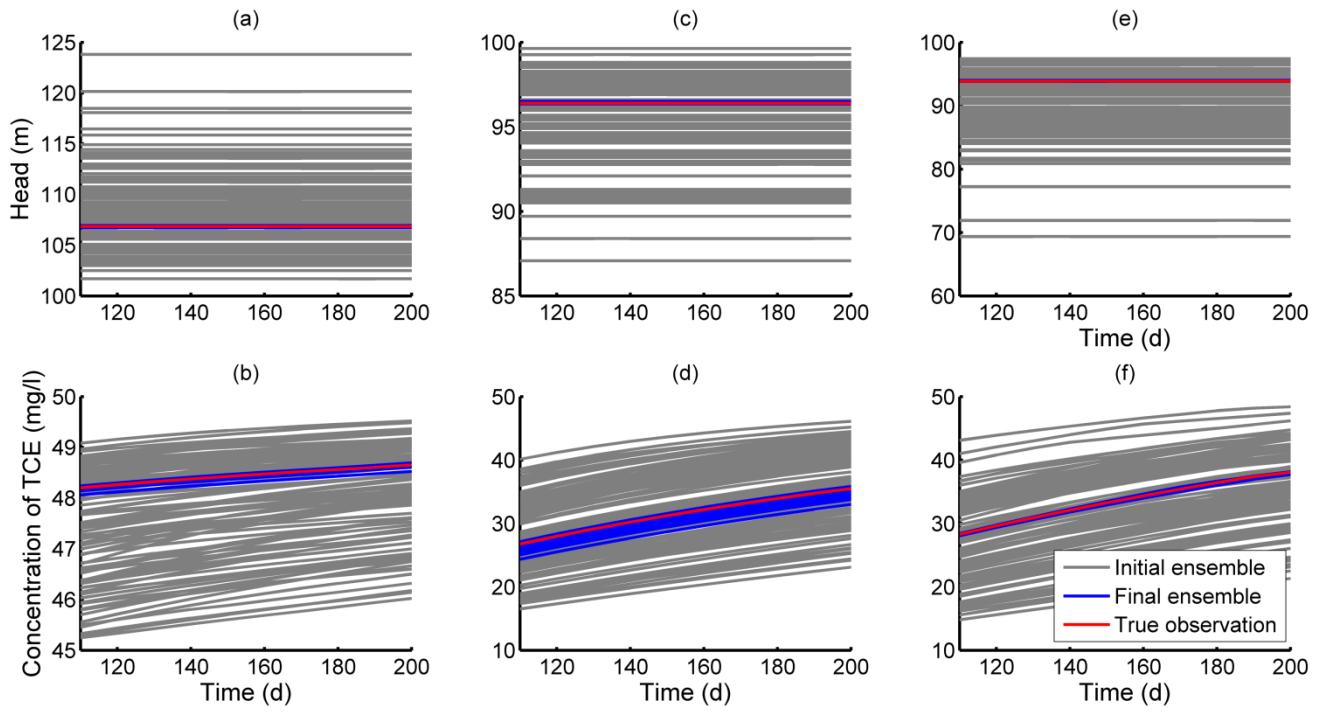
2

3

4



1 **Fig. 11** The performance of data match. (a), (c) and (e) show the data match of the head of three selected
2 wells respectively, while (b), (d) and (f) show the data match of the TCE concentration data of three
3 selected wells respectively.
4
5



1 **Fig. 12** The performance of model prediction. (a), (c) and (e) show the prediction of the head of three
2 selected wells respectively, while (b), (d) and (f) show the prediction of the TCE concentration data of
3 three selected wells respectively.
4
5

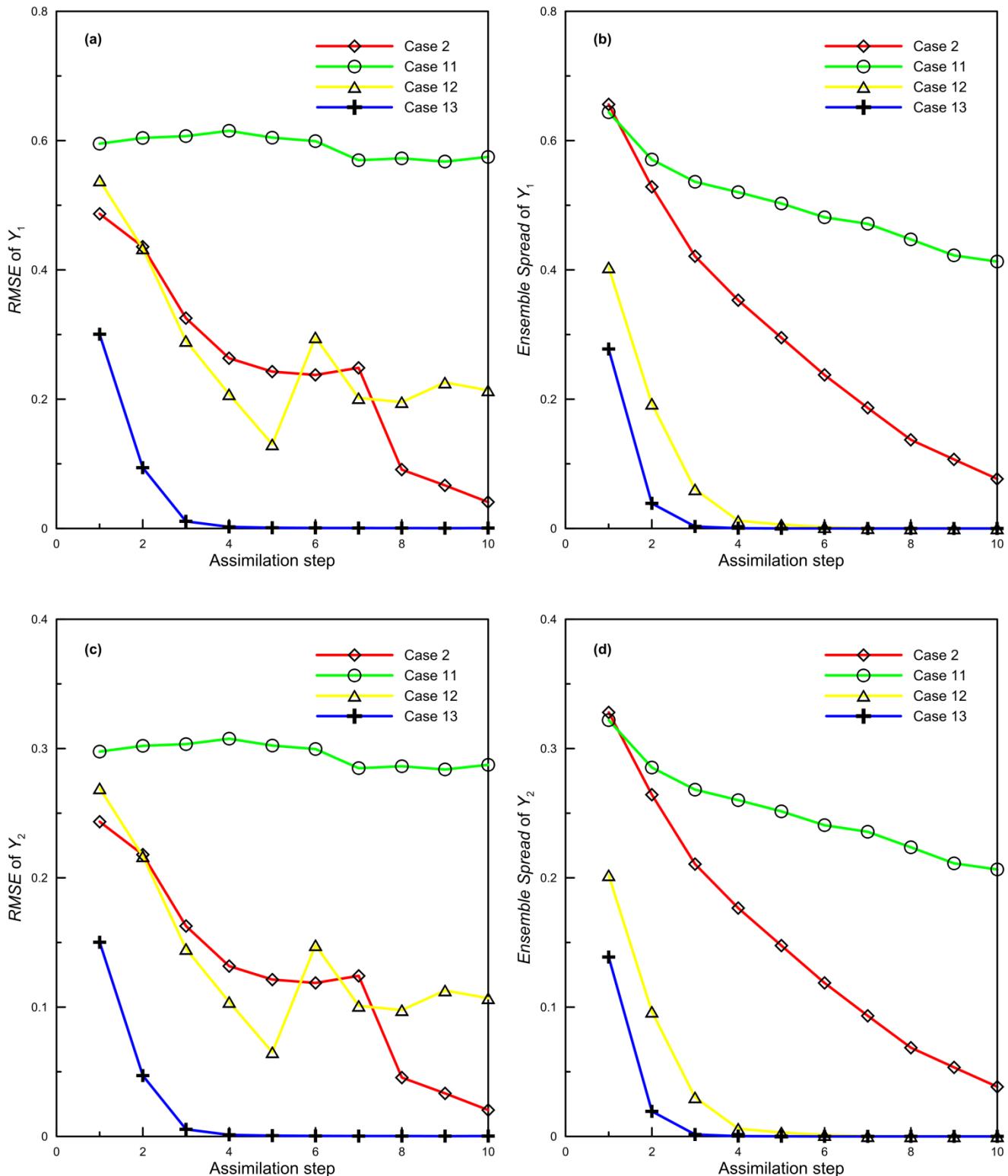


Fig. 13 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different sampling strategies.

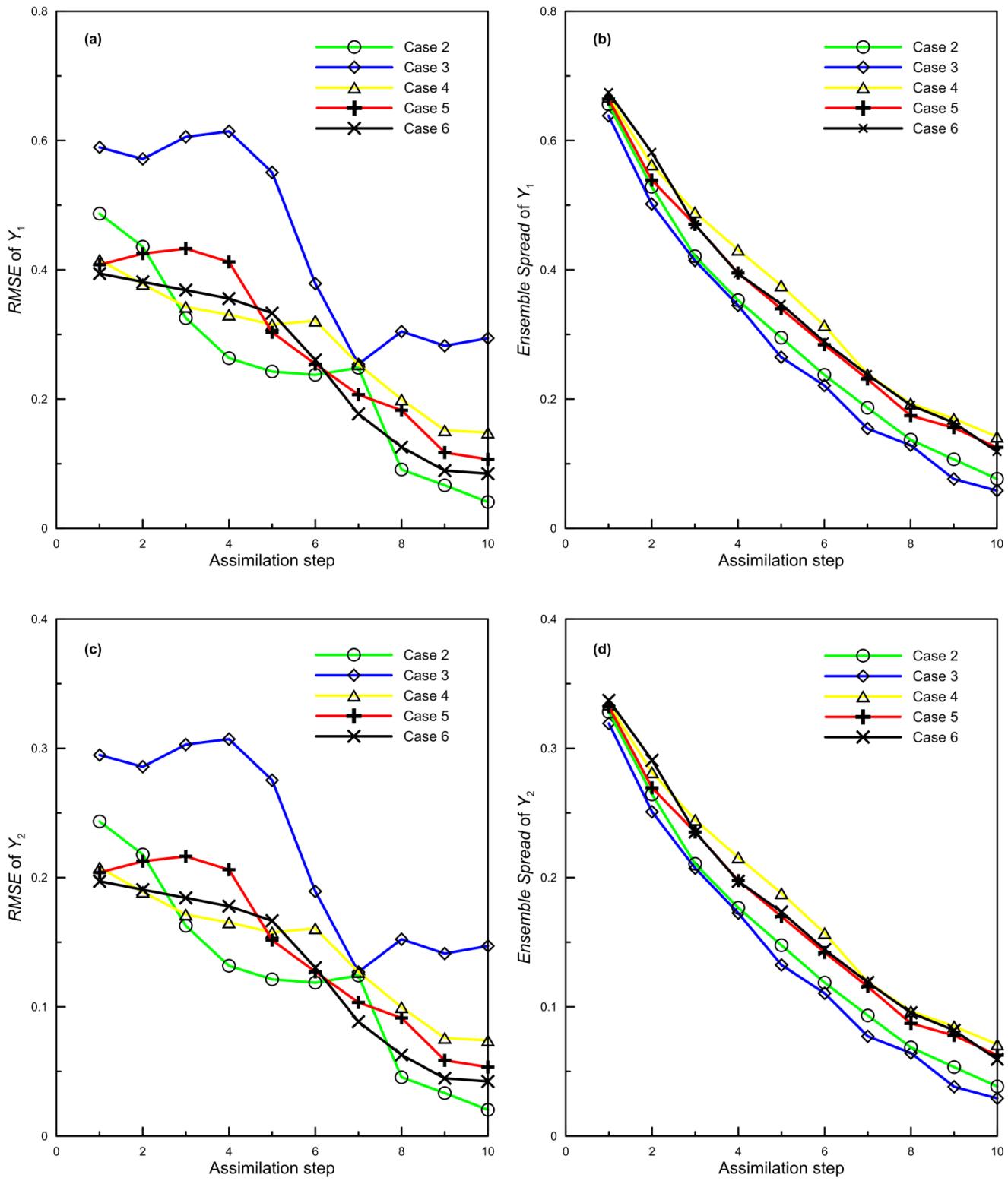


Fig. 13-14 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different ensemble sizes.

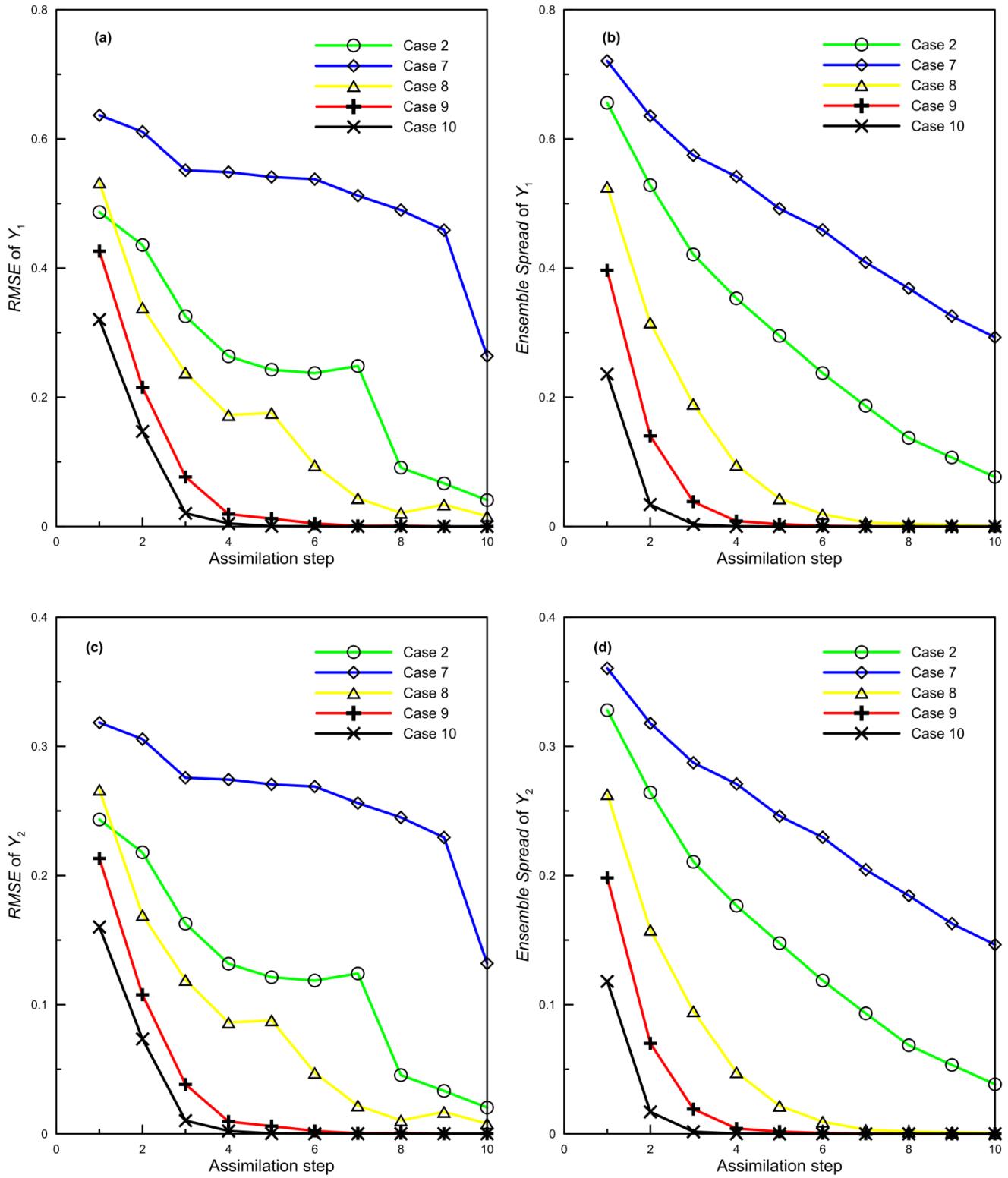


Fig. 14–15 The performance indicators (the RMSE and the Ensemble Spread) at each assimilation step for different numbers of optimal sampling locations.

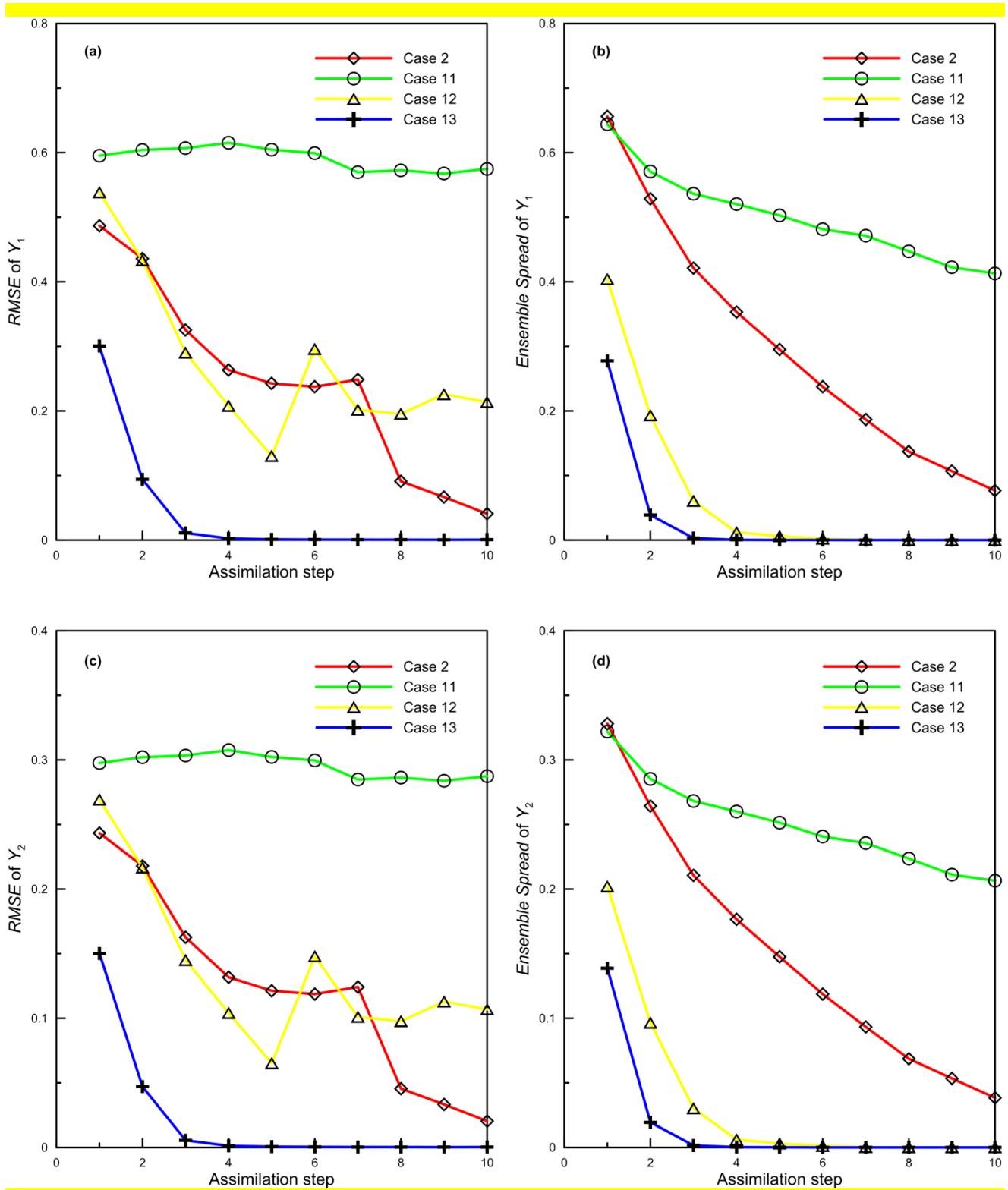


Fig. 15 The performance indicators (RMSE and Ensemble Spread) at each assimilation step for different sampling strategies.

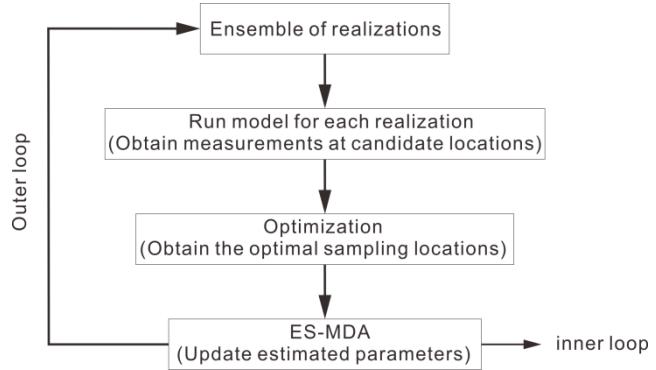


Fig. 16 The loop diagram of the improved SEOD method.

1

2

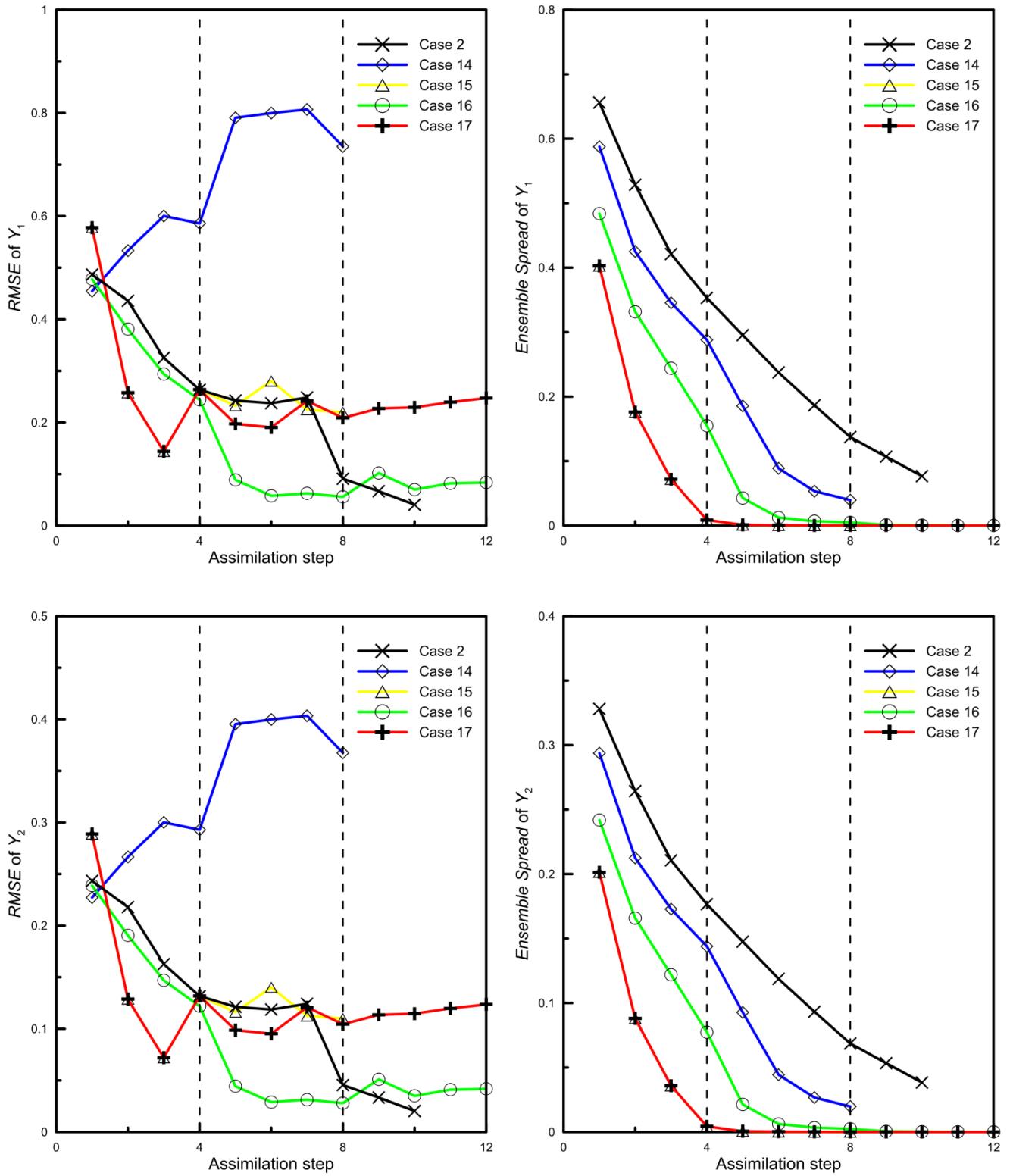
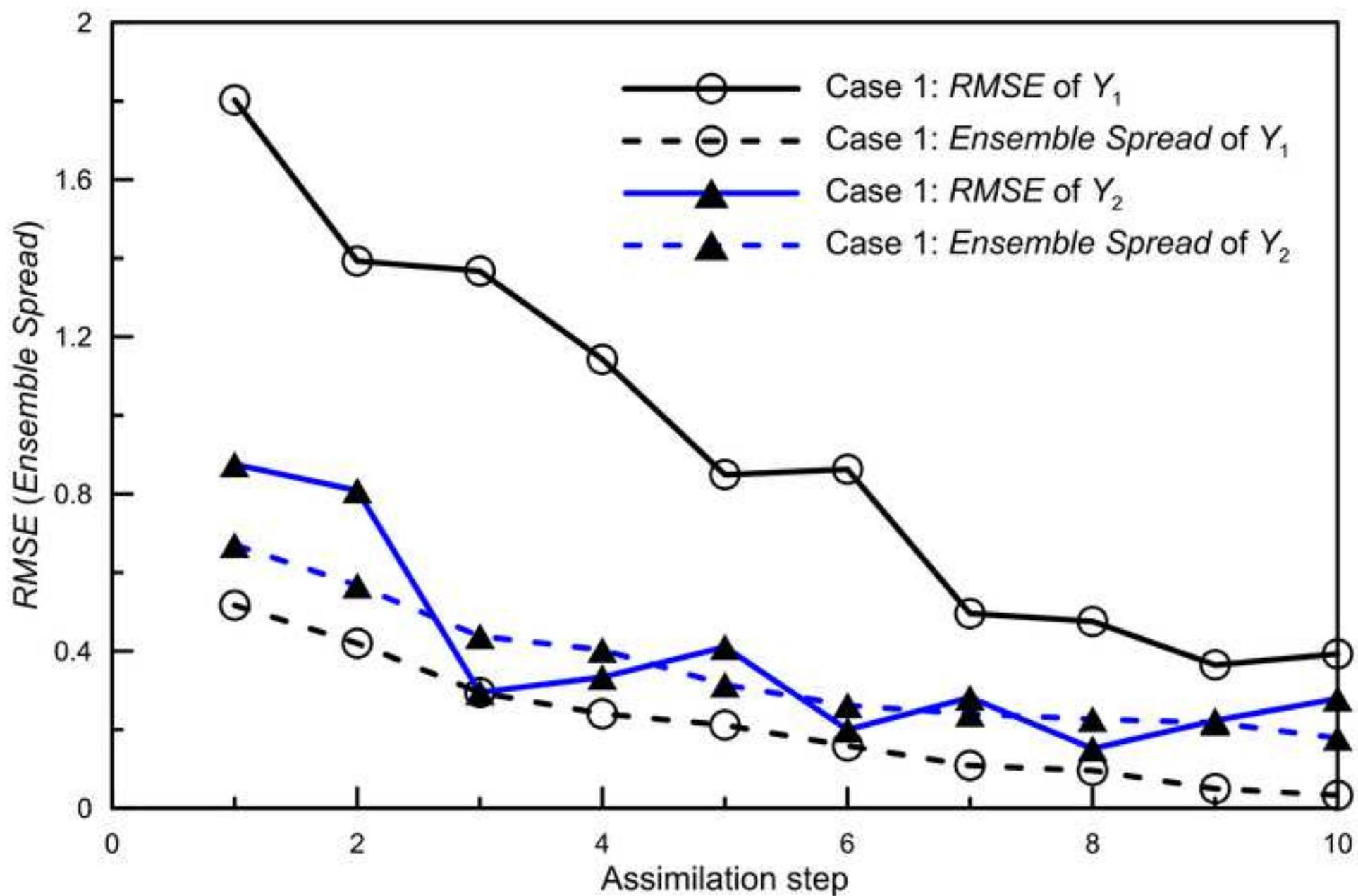
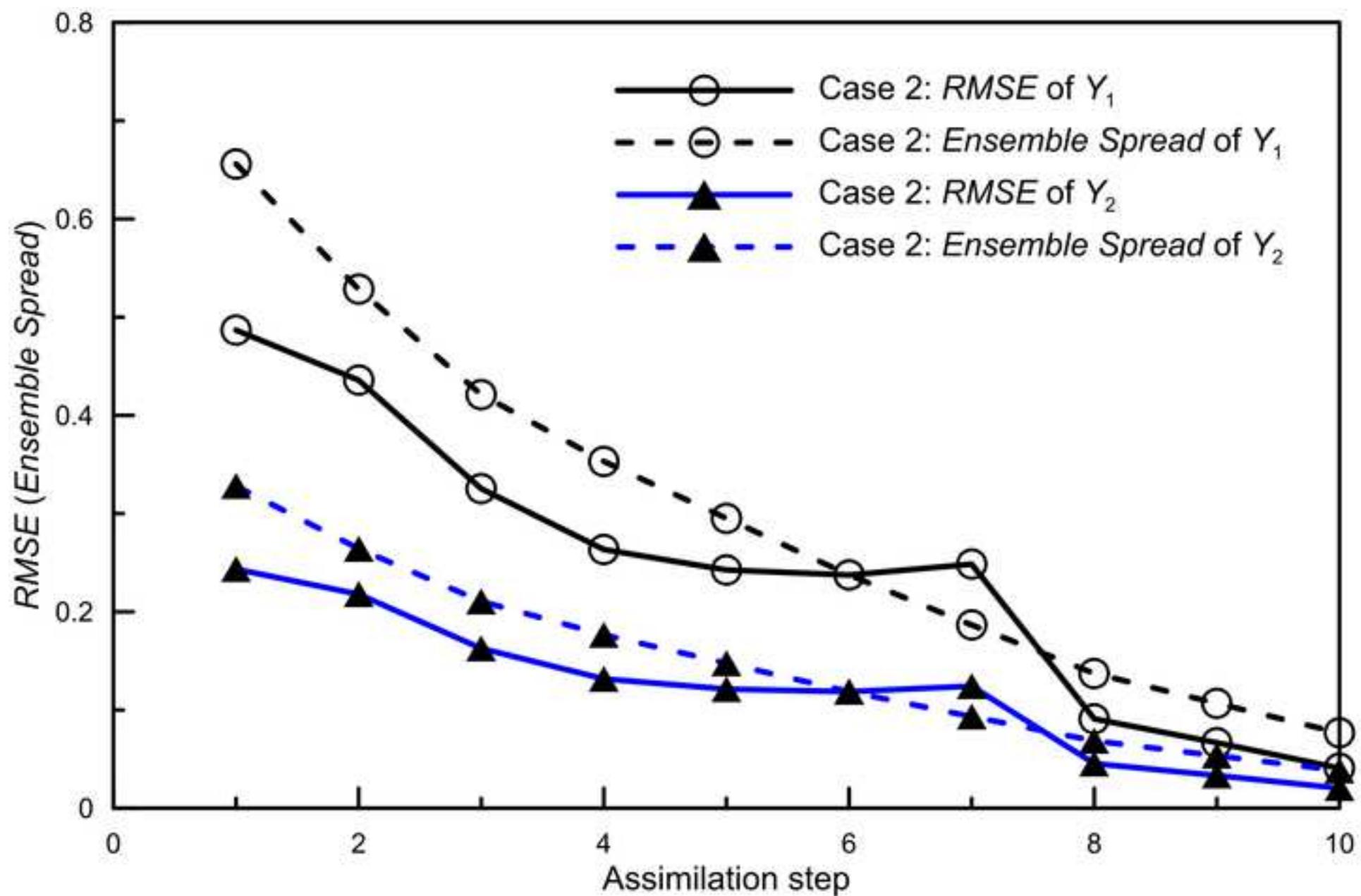
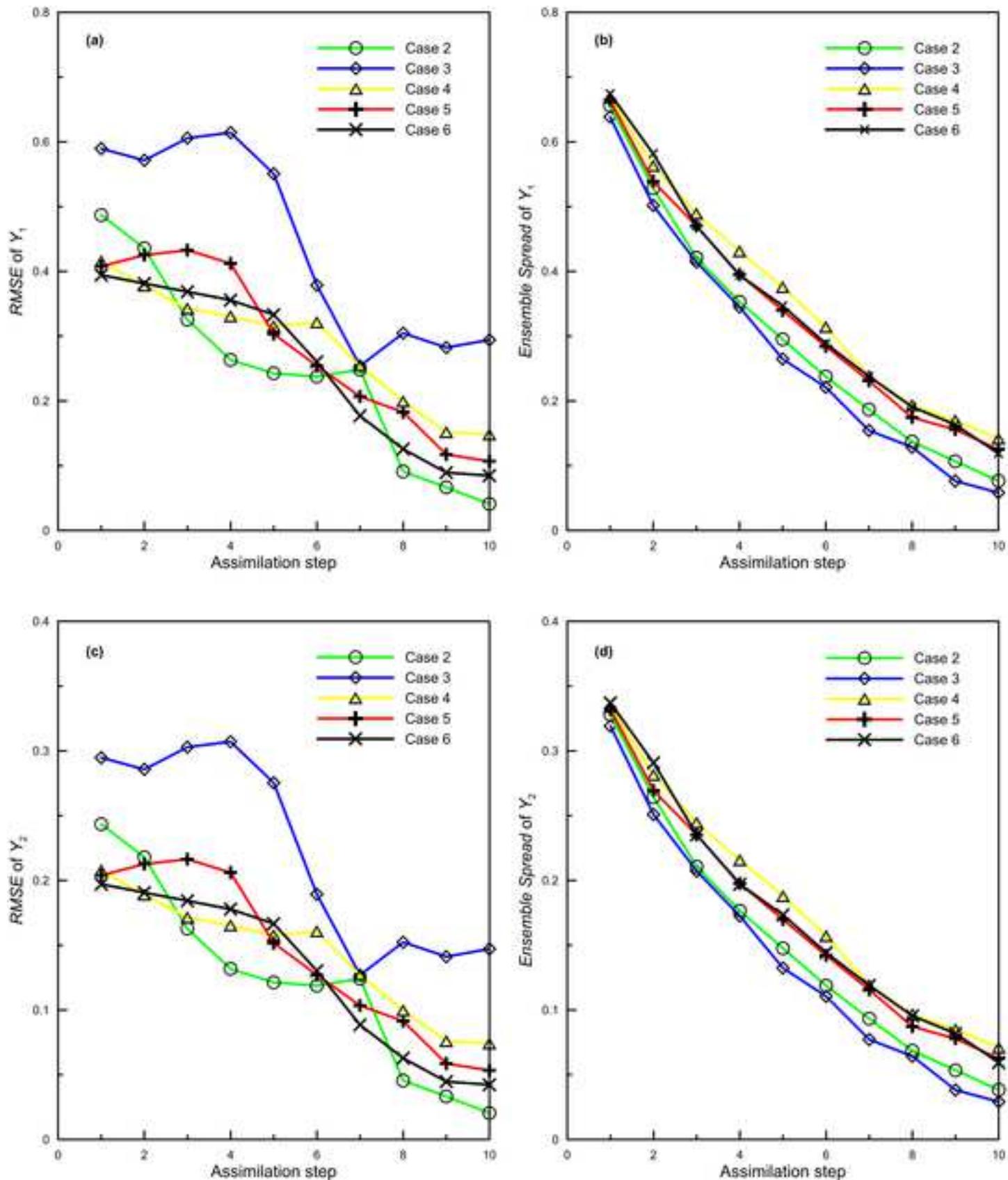
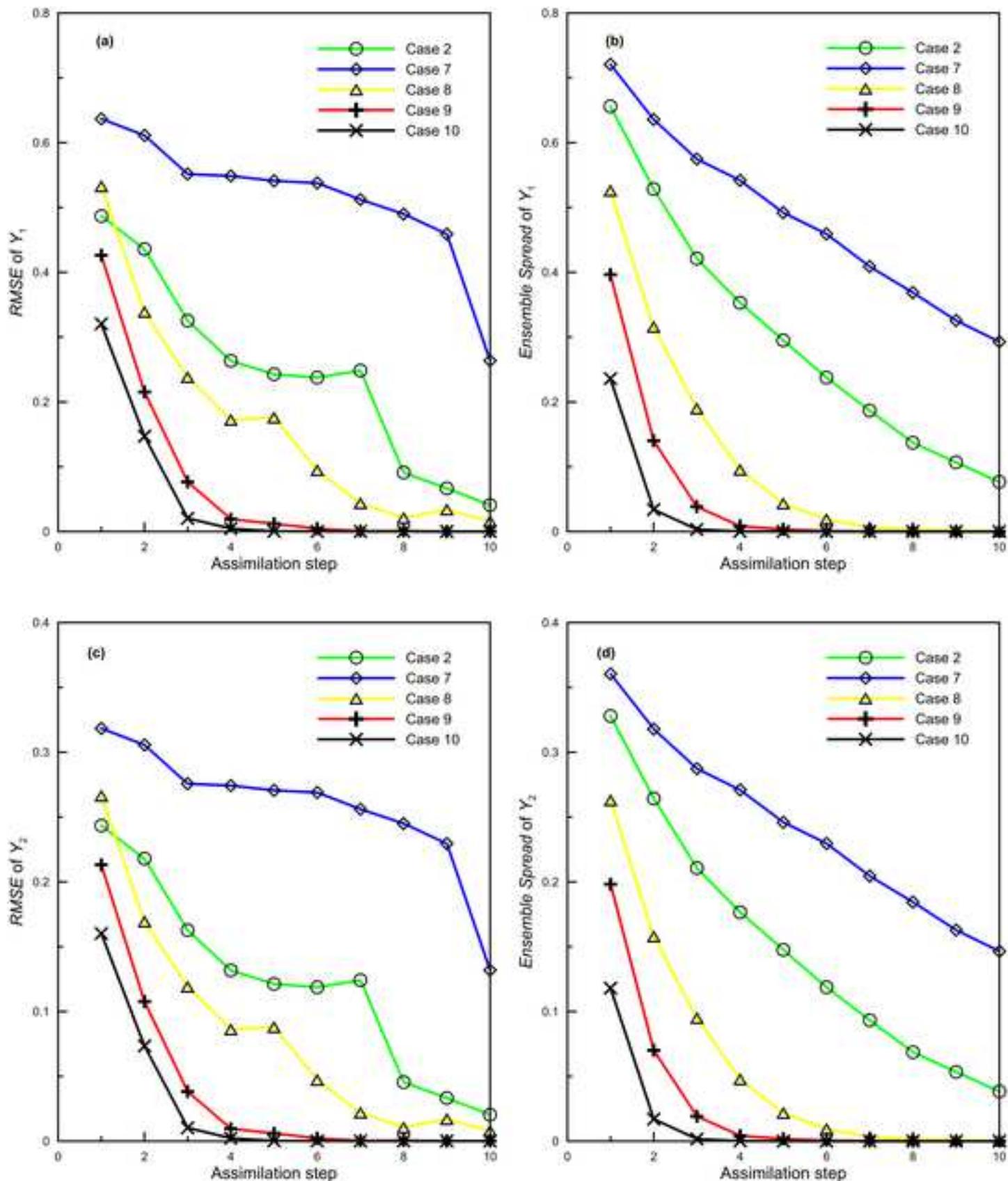


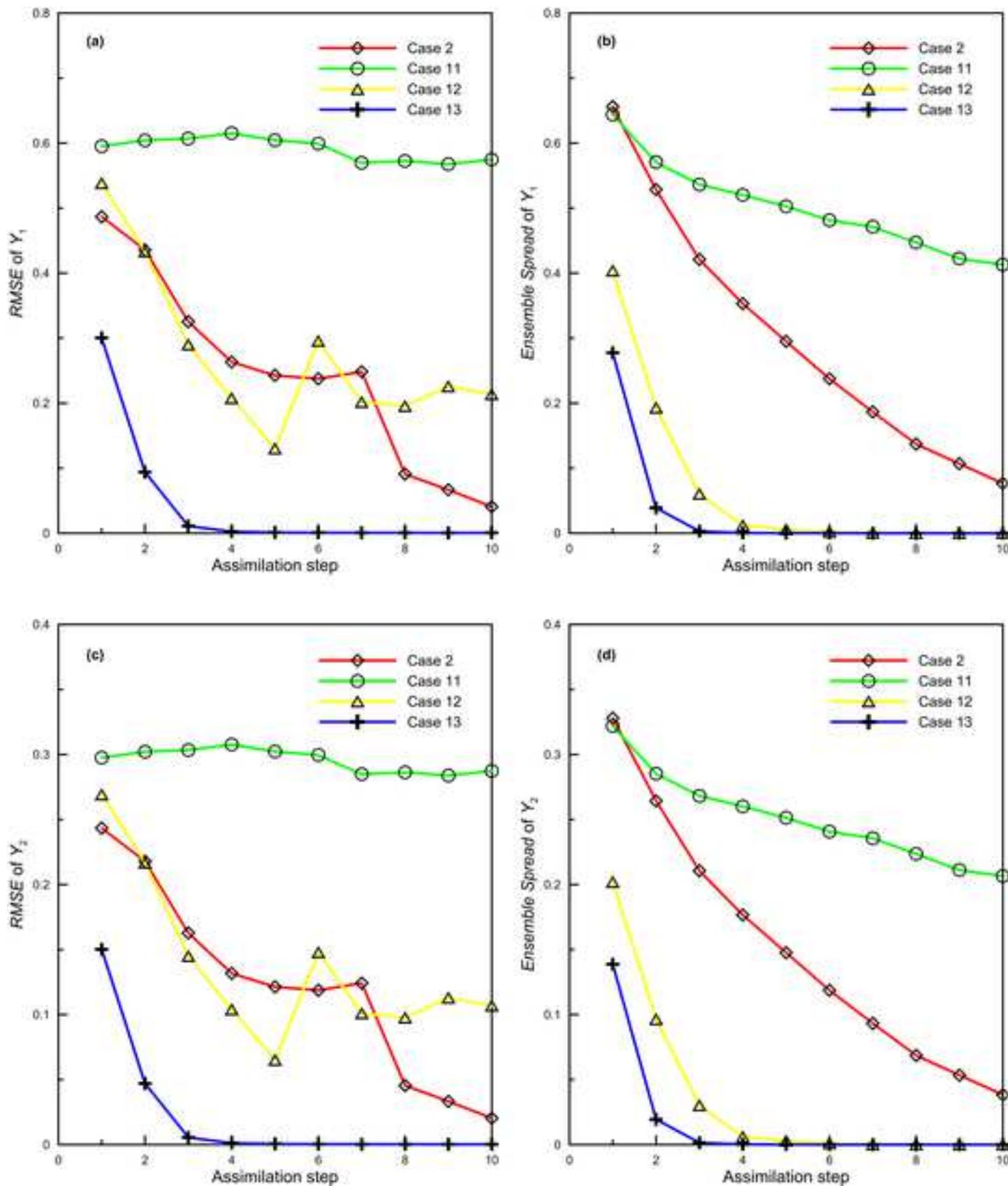
Fig. 17 The performance indicators (the *RMSE* and the *Ensemble Spread*) at each assimilation step for different cases using the improved SEOD method.

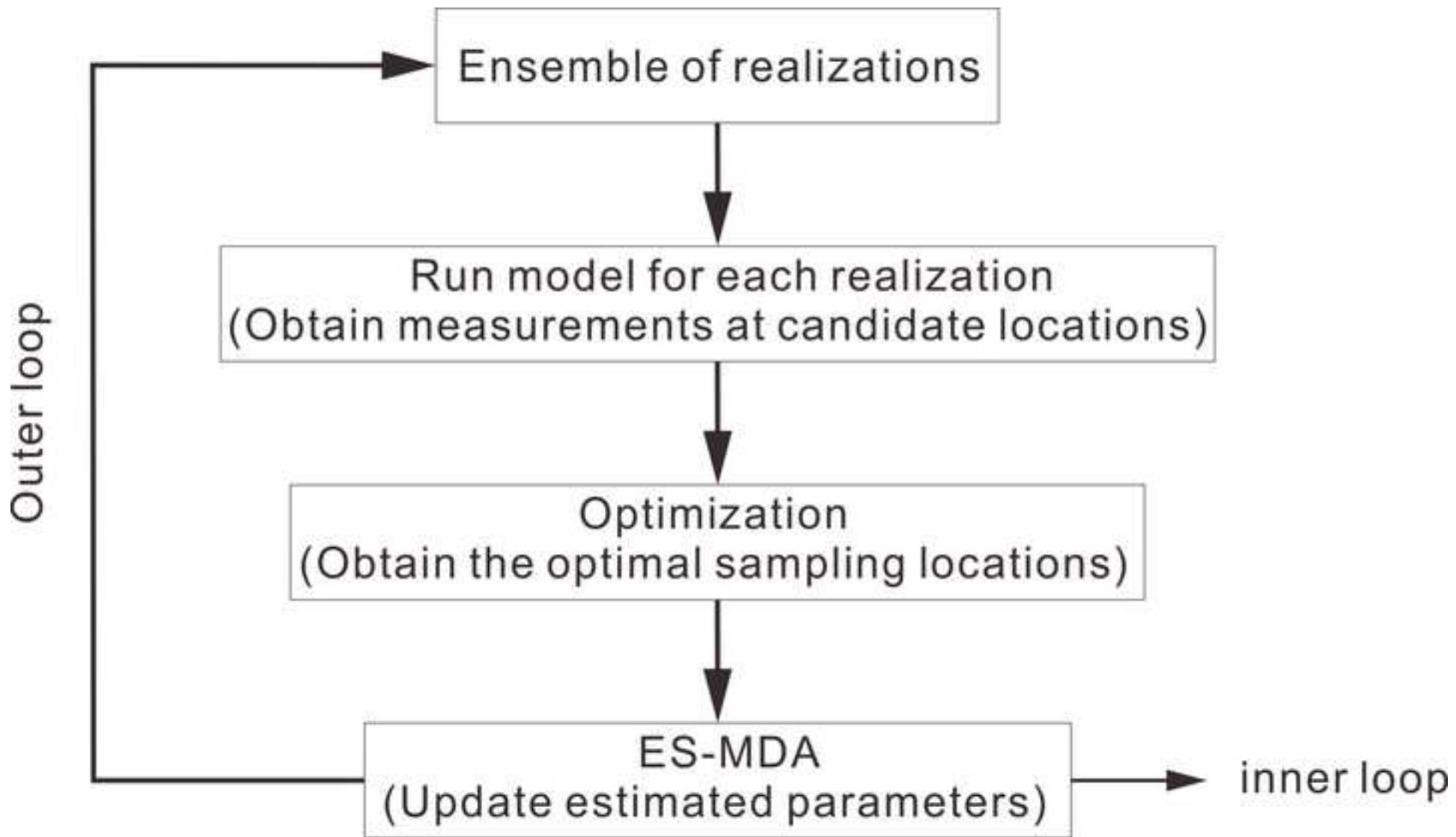


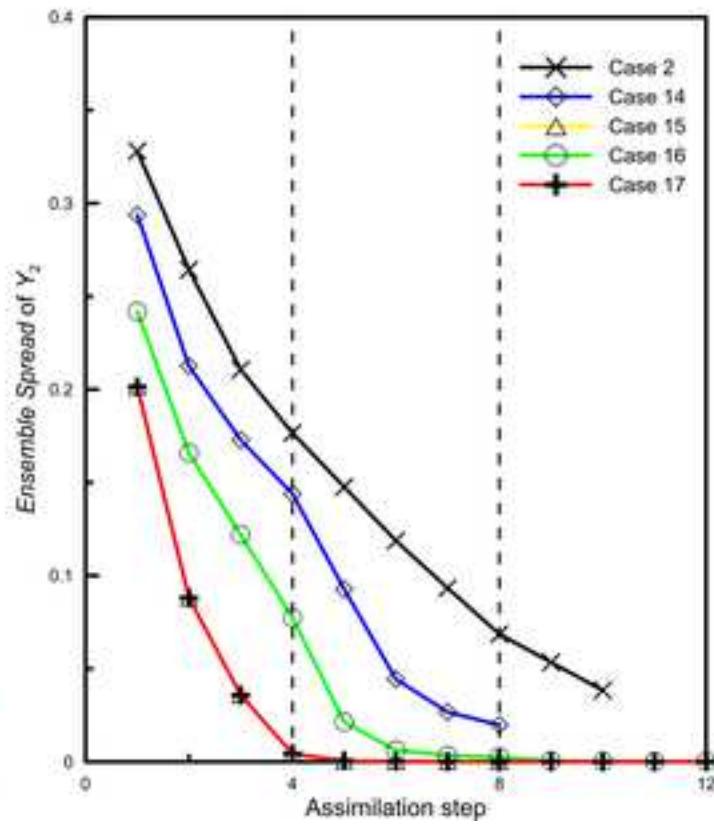
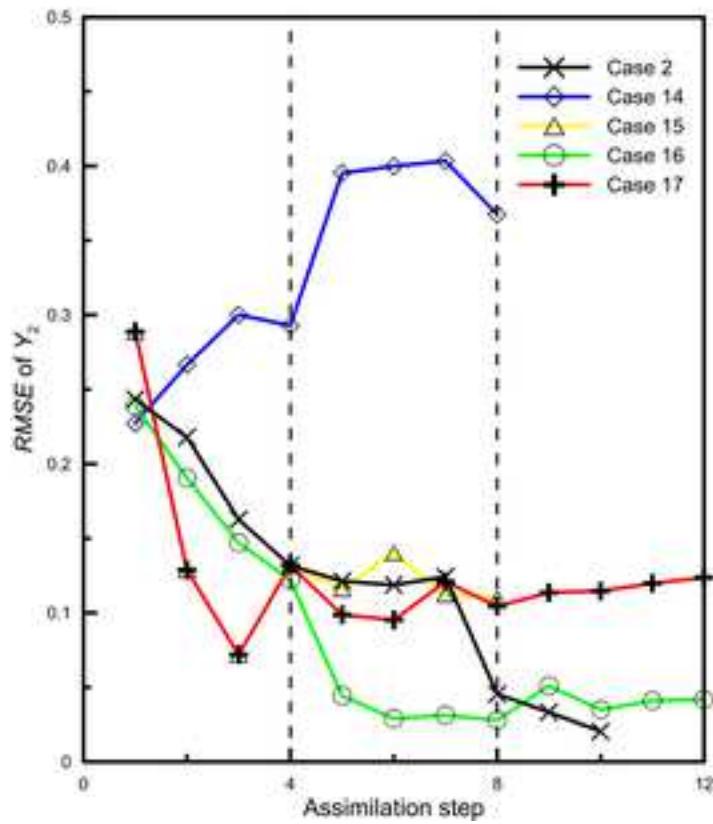
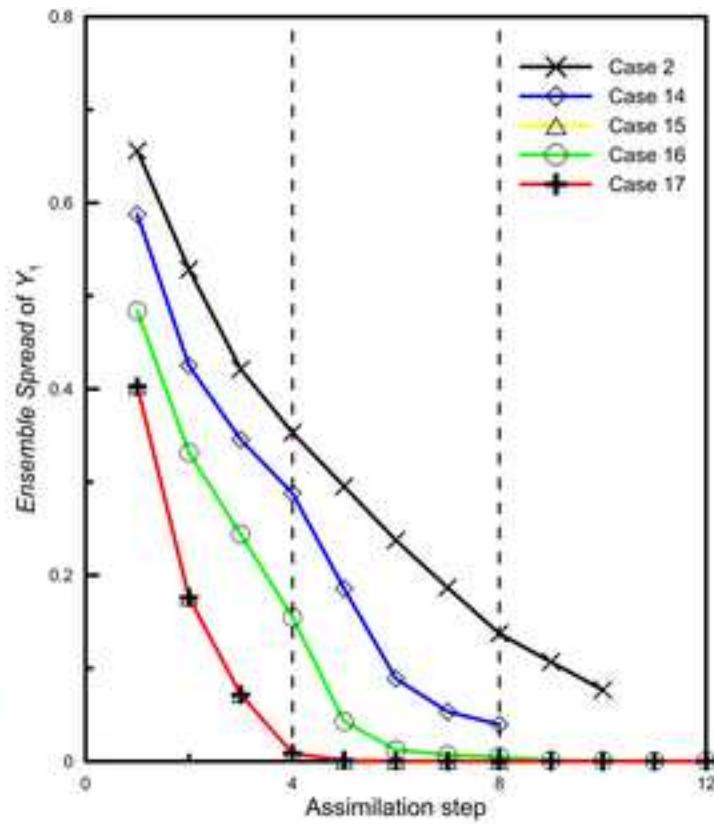
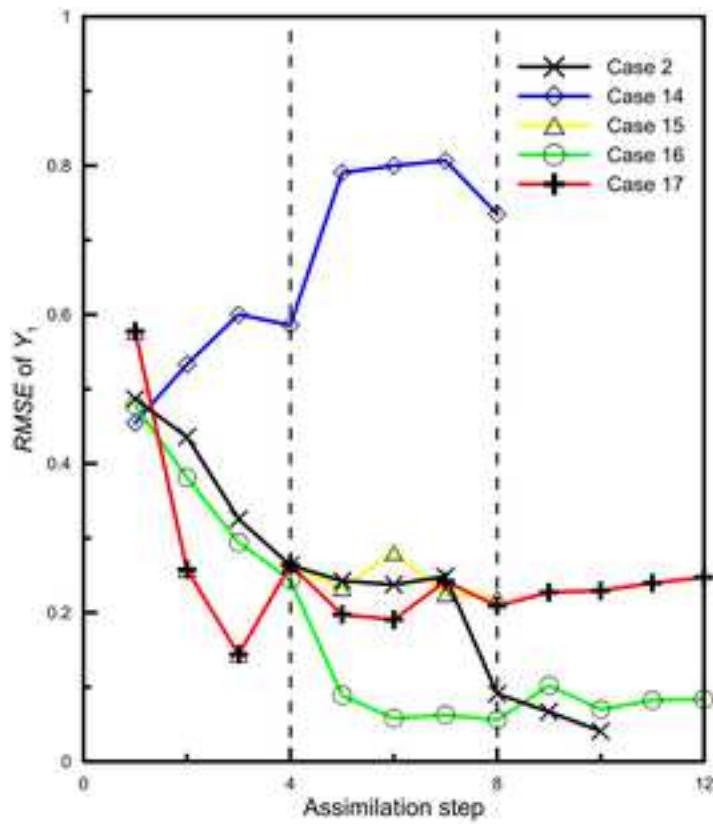


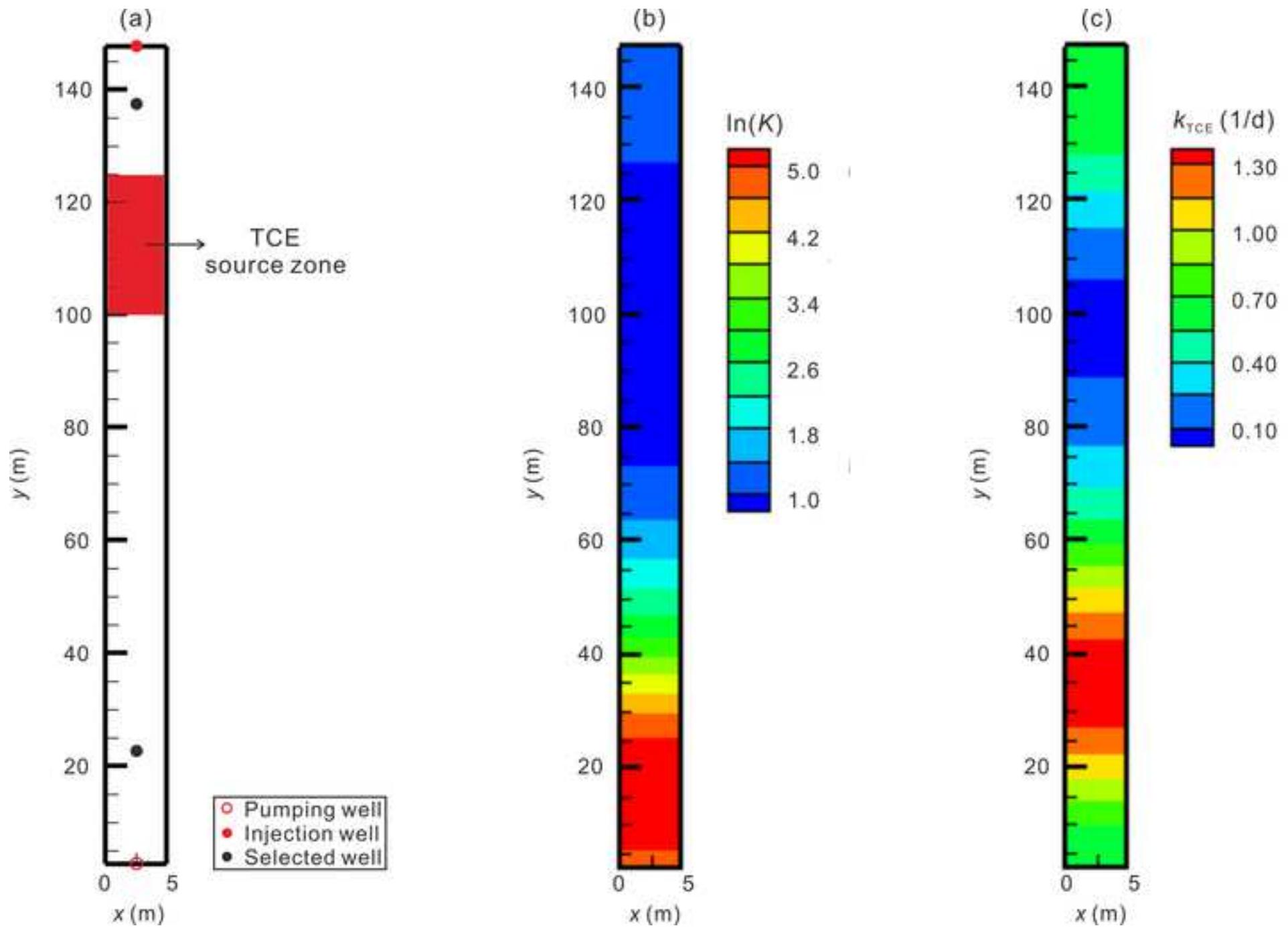


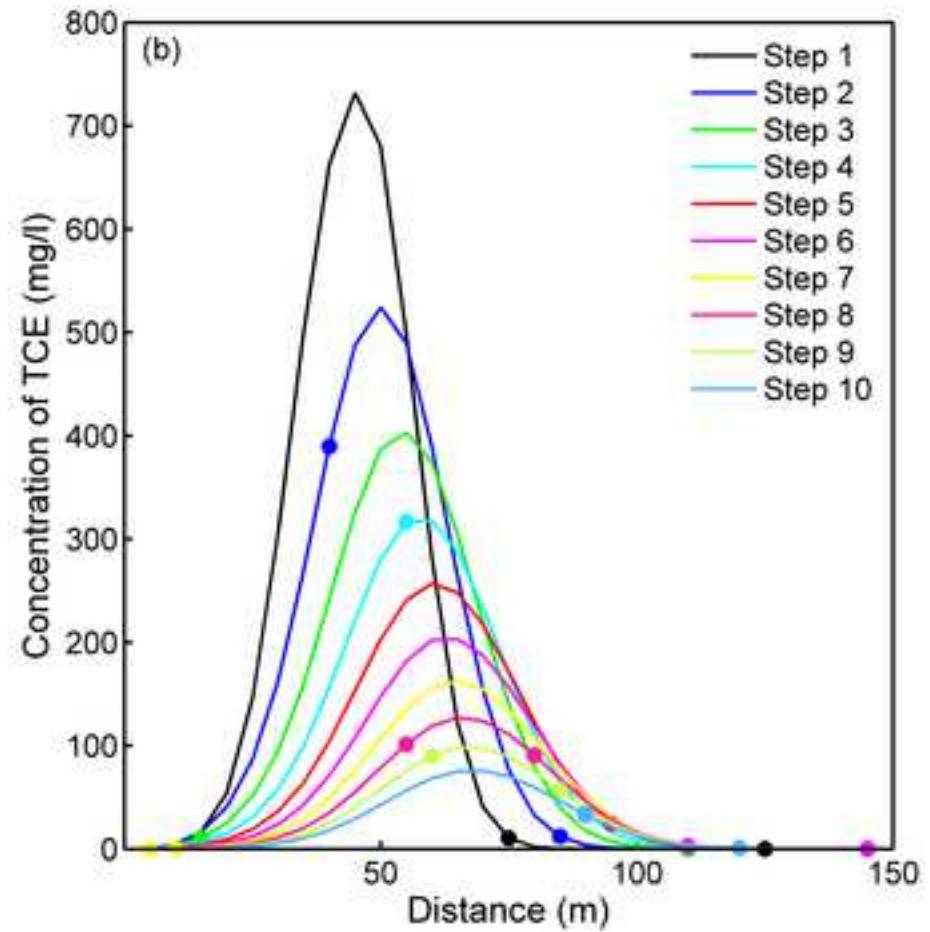
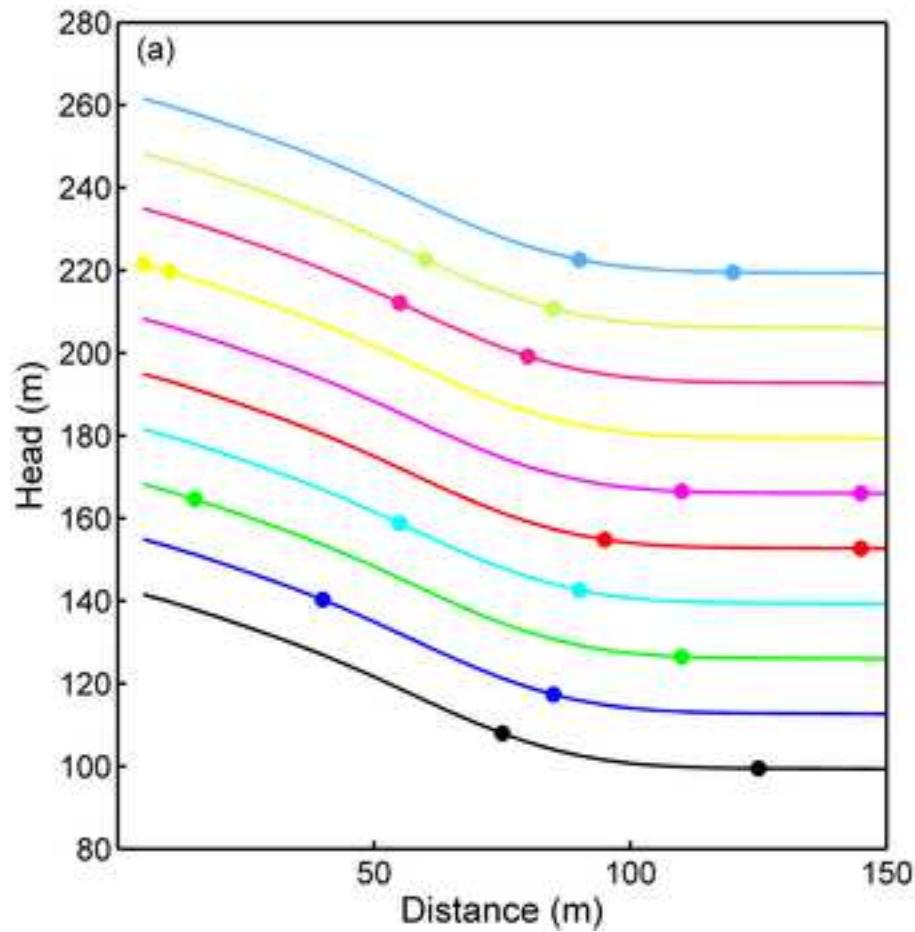


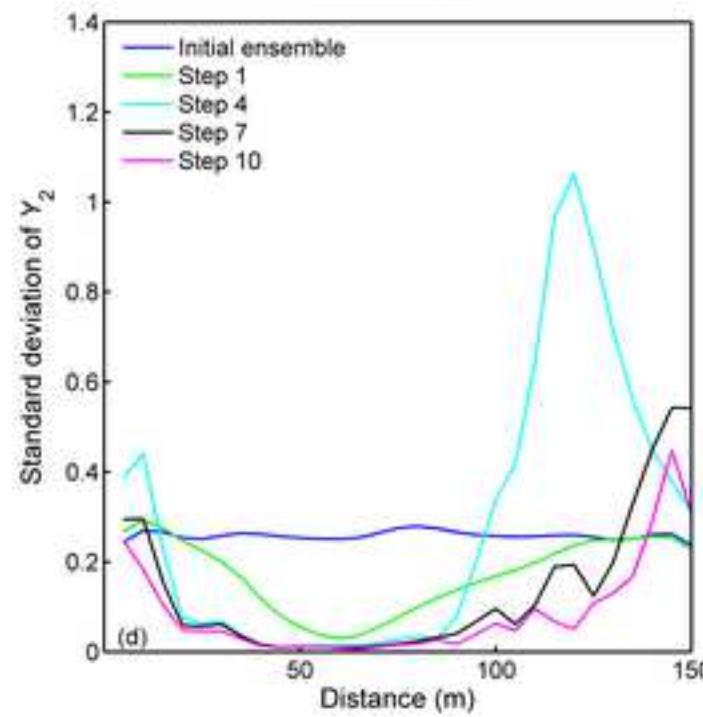
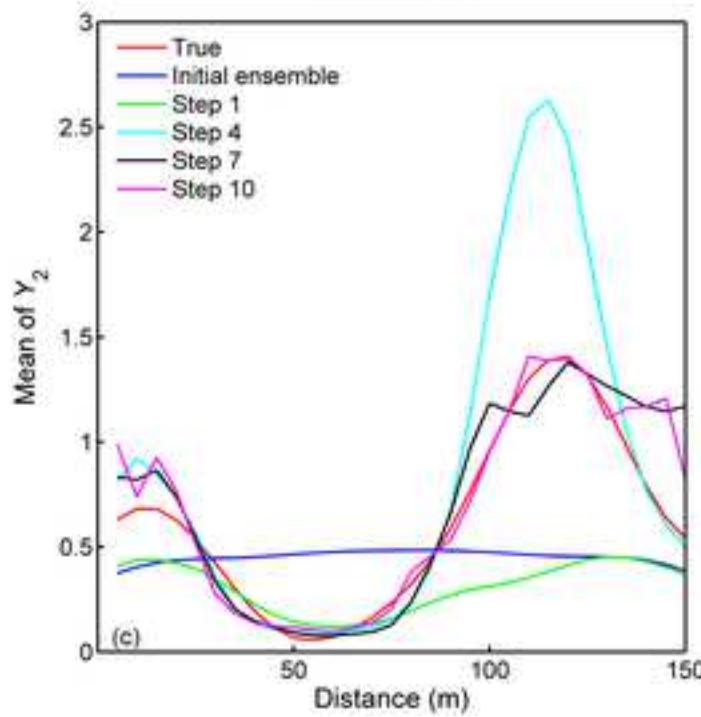
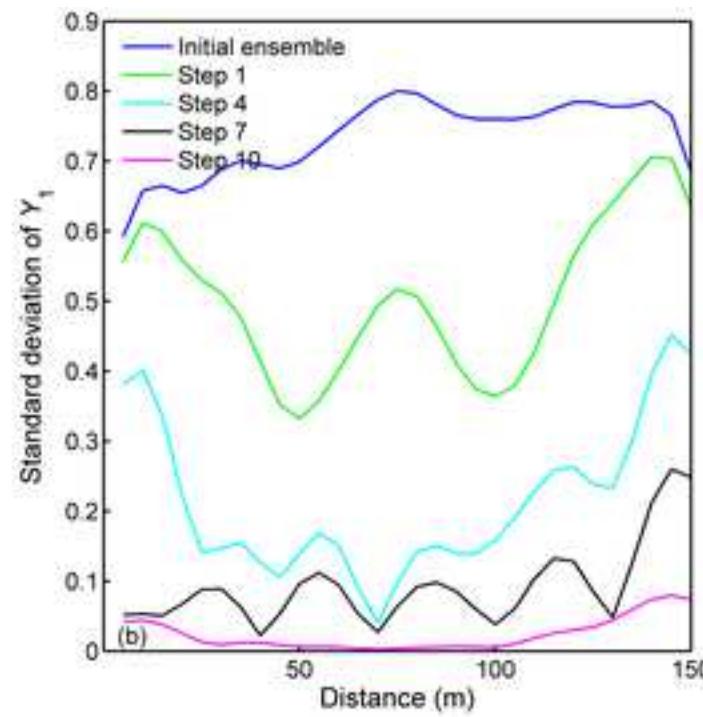
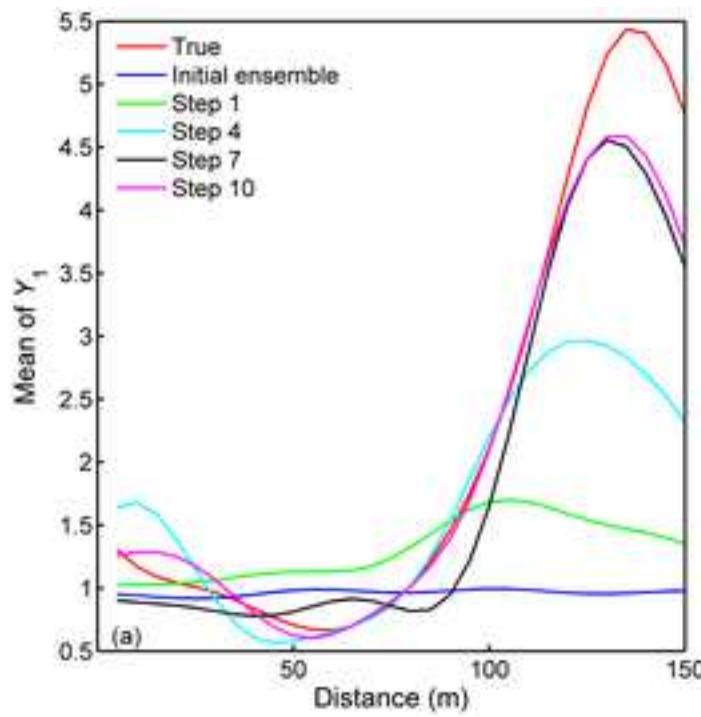


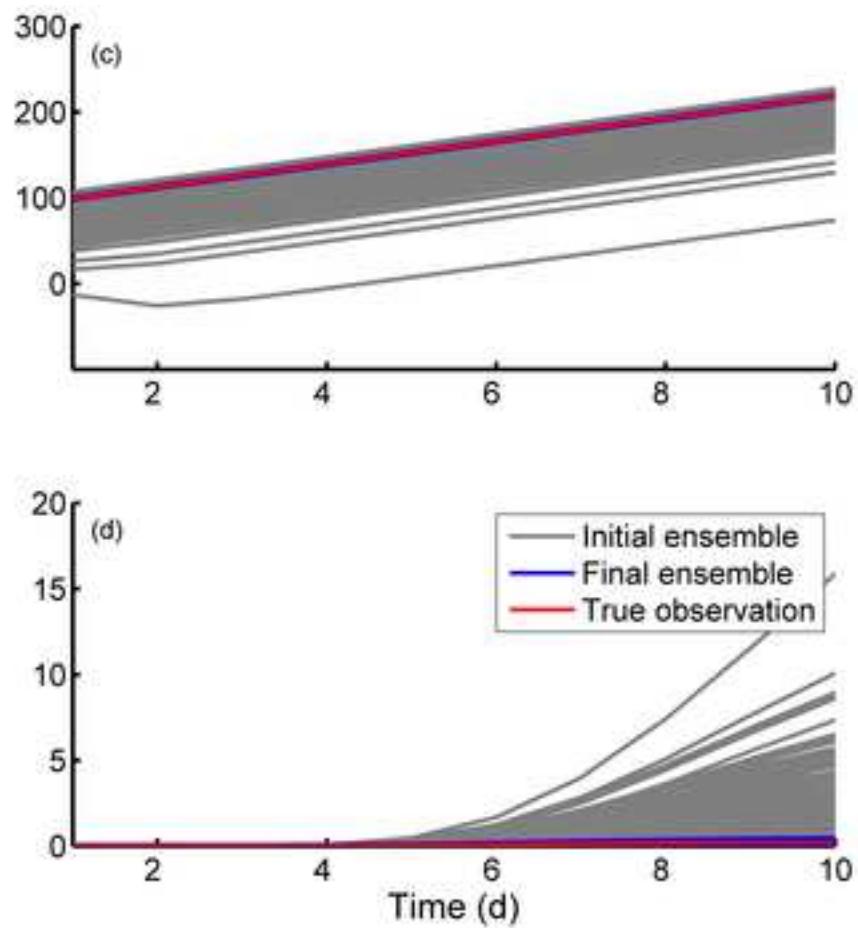
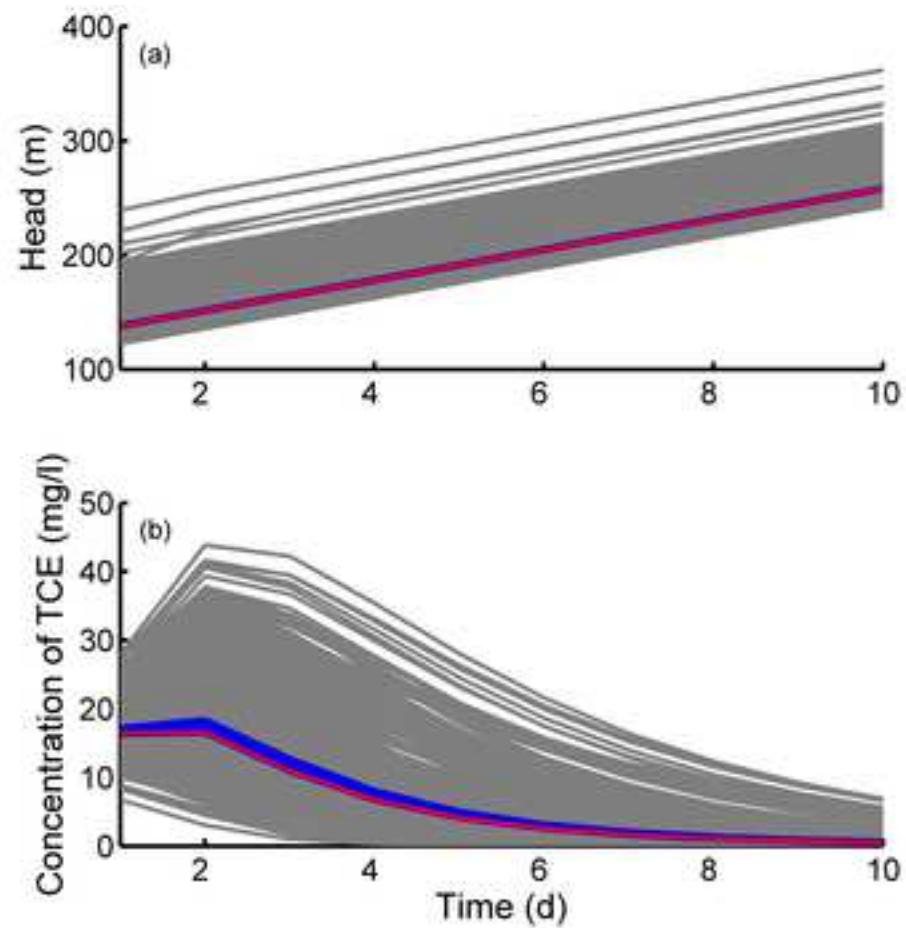


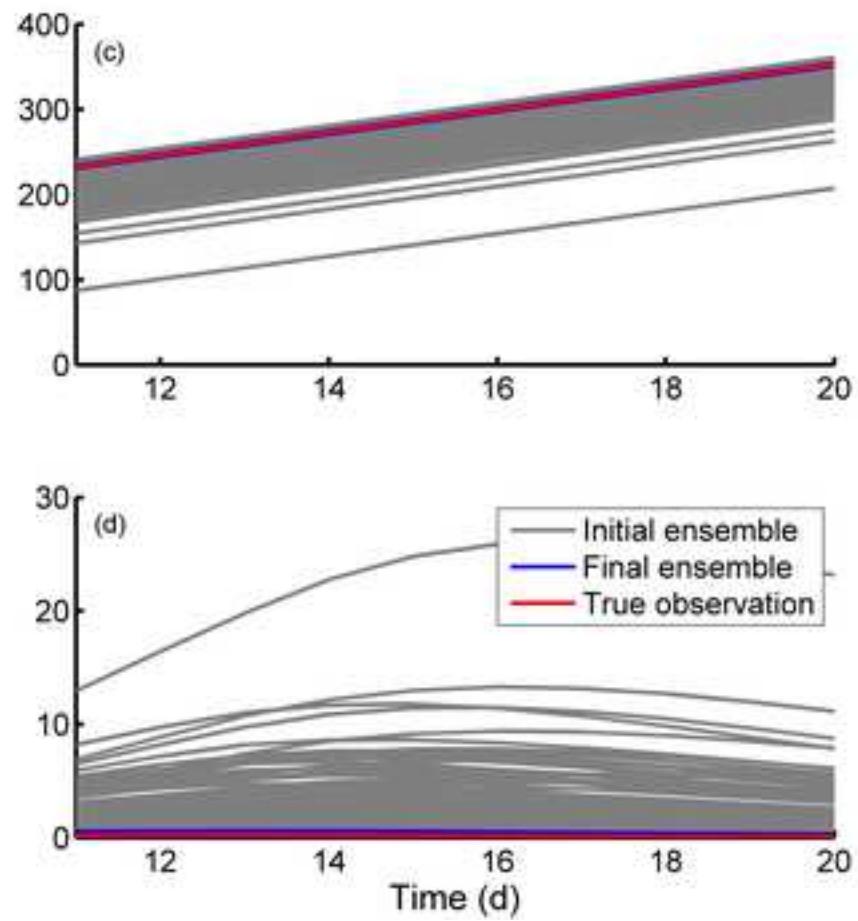
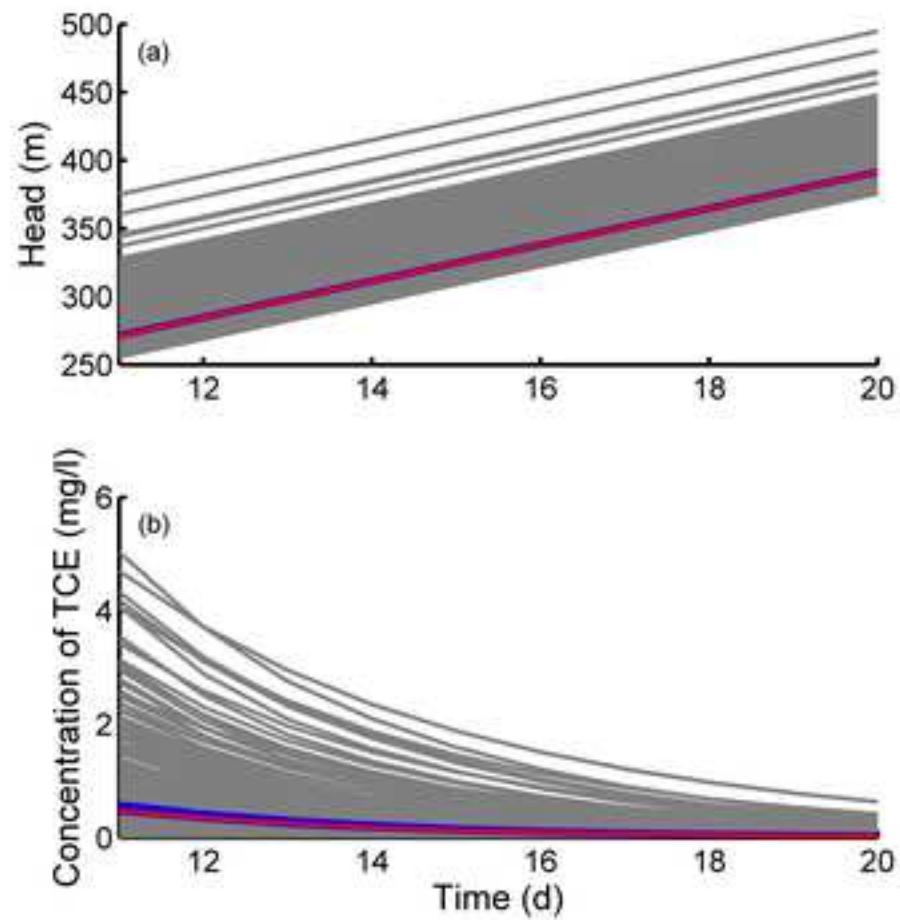


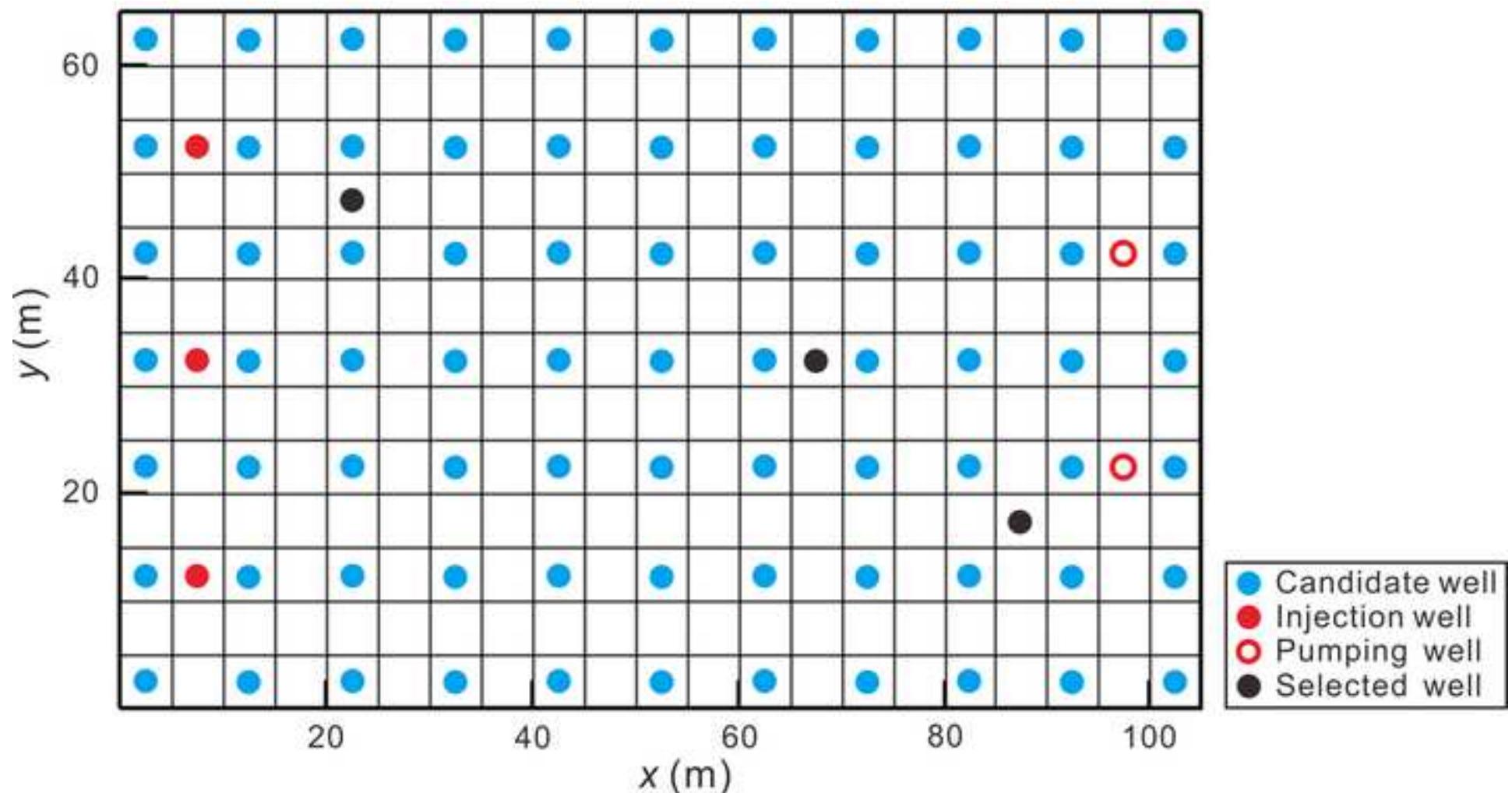


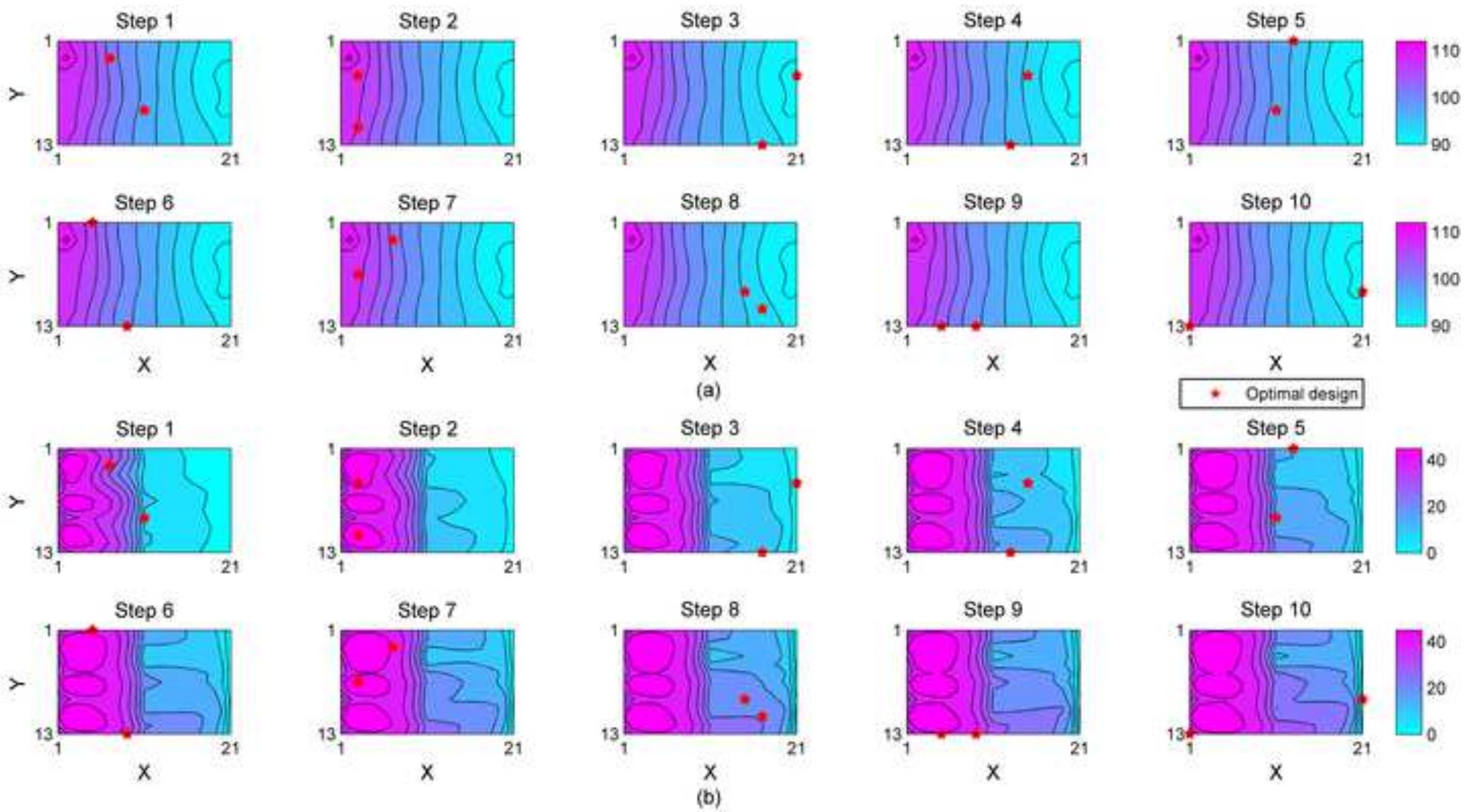


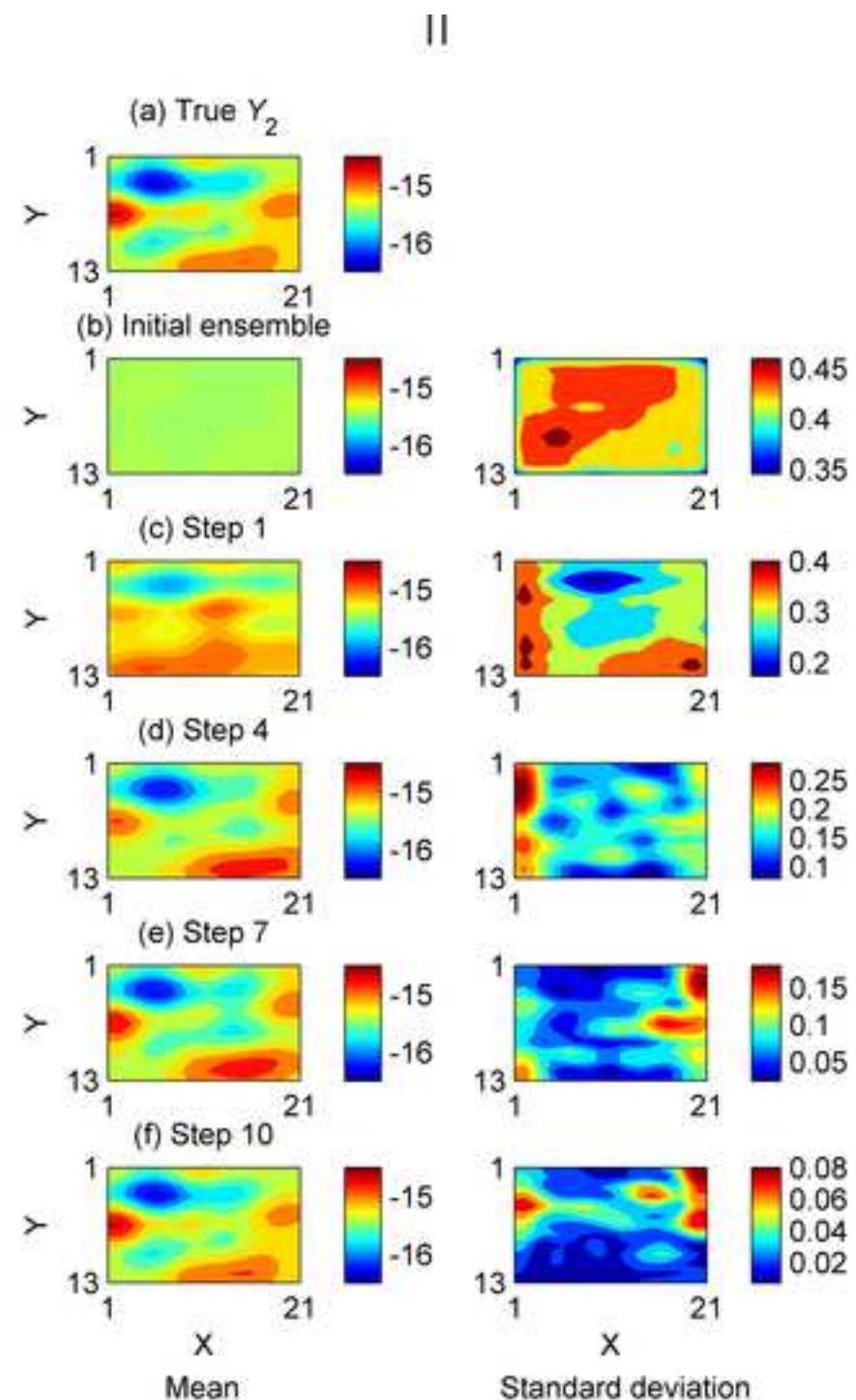
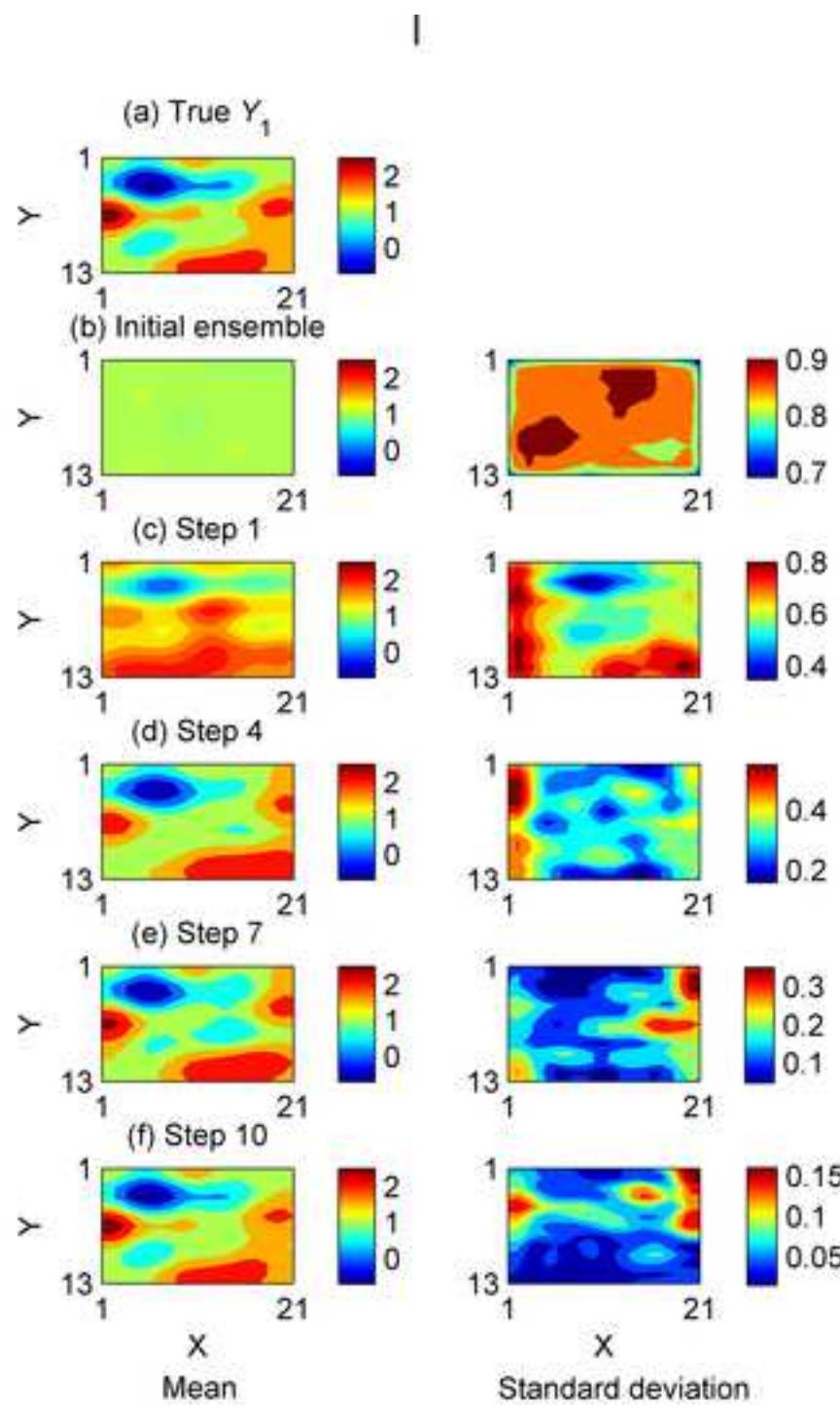


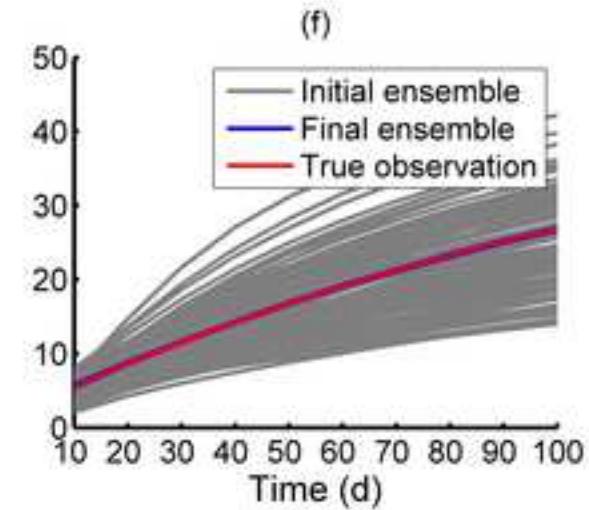
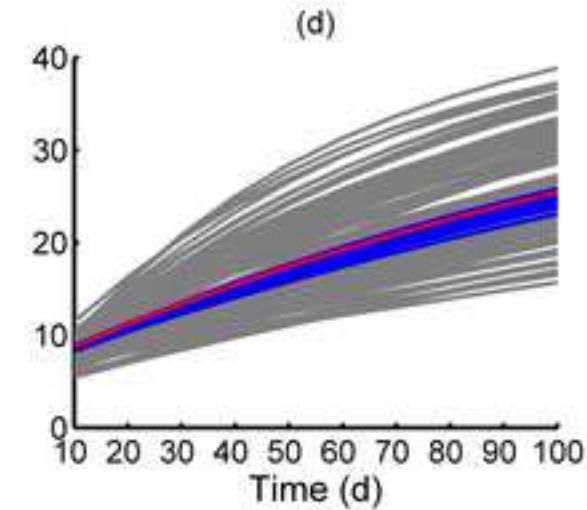
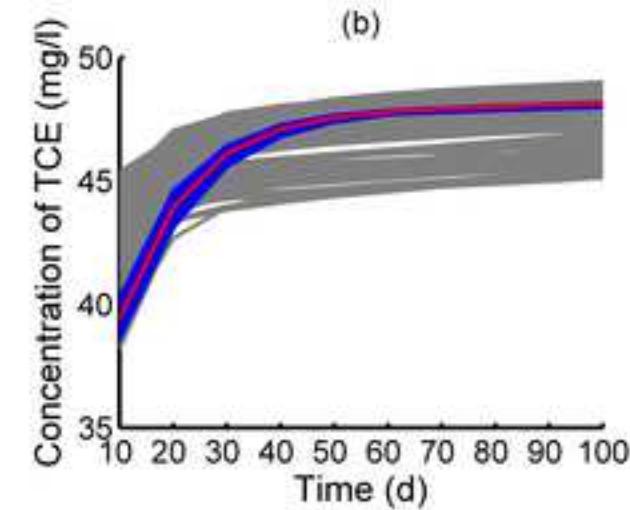
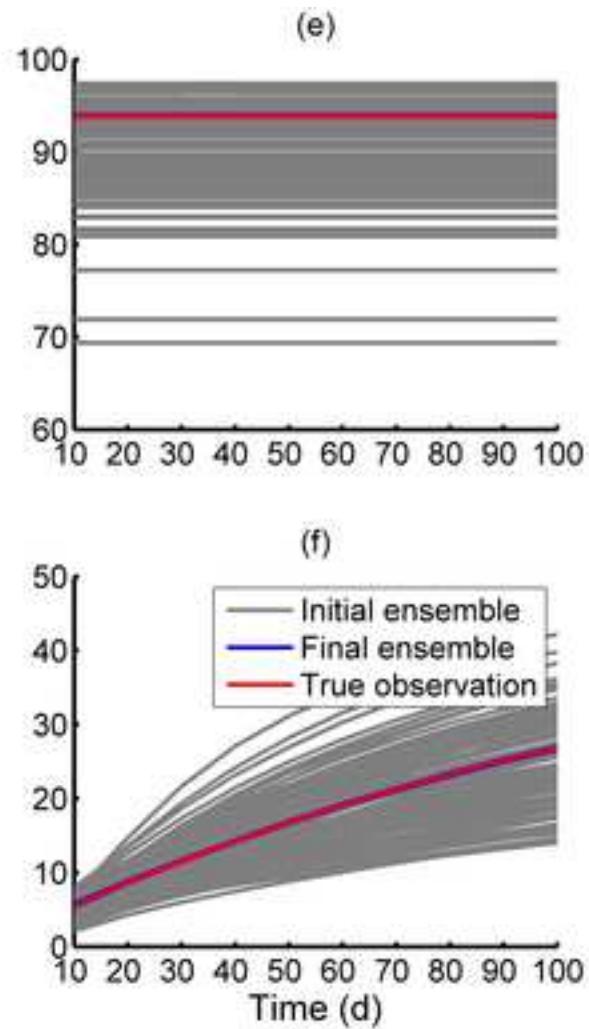
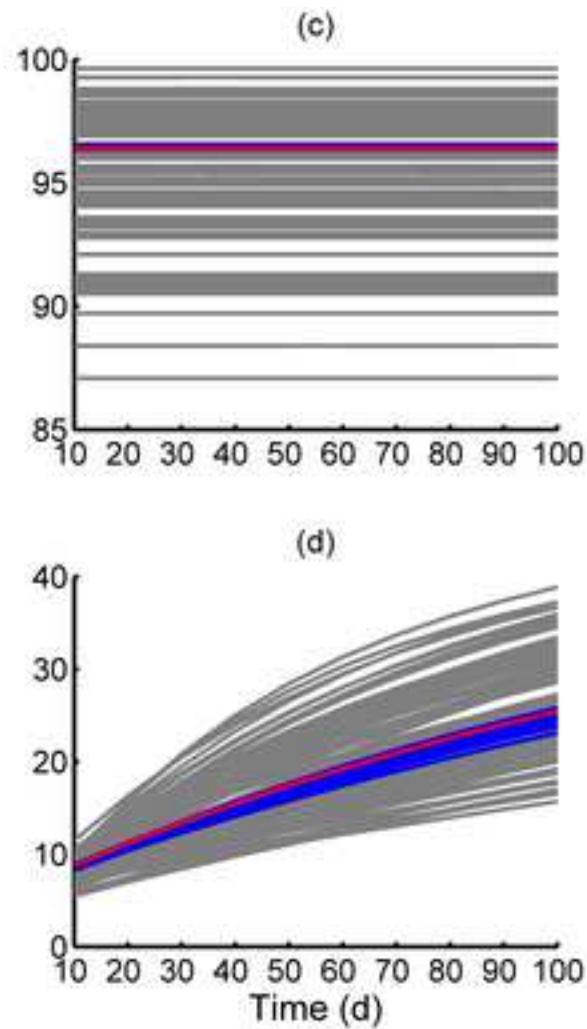
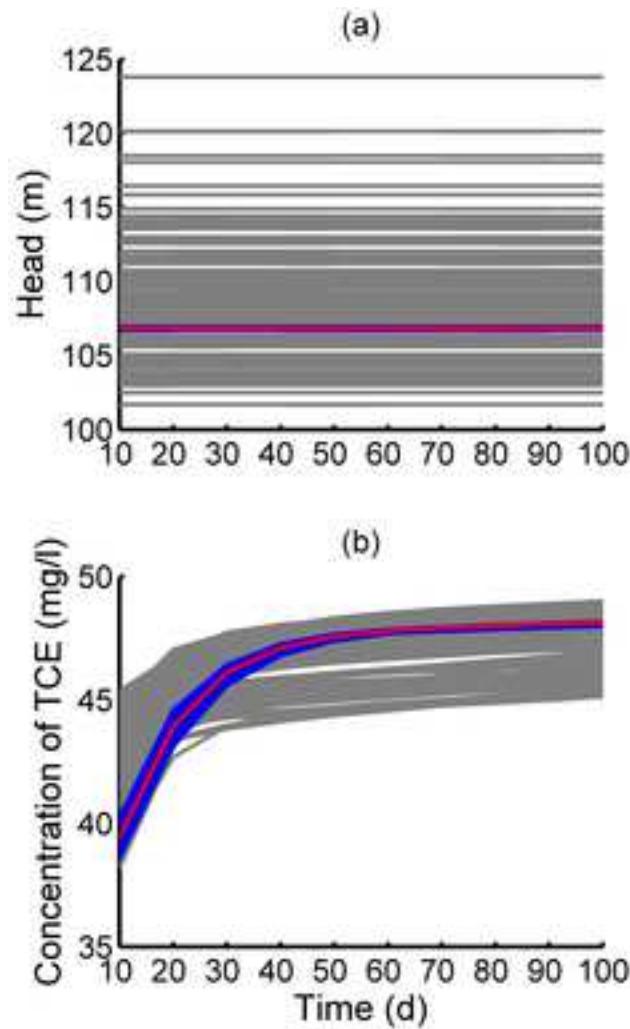


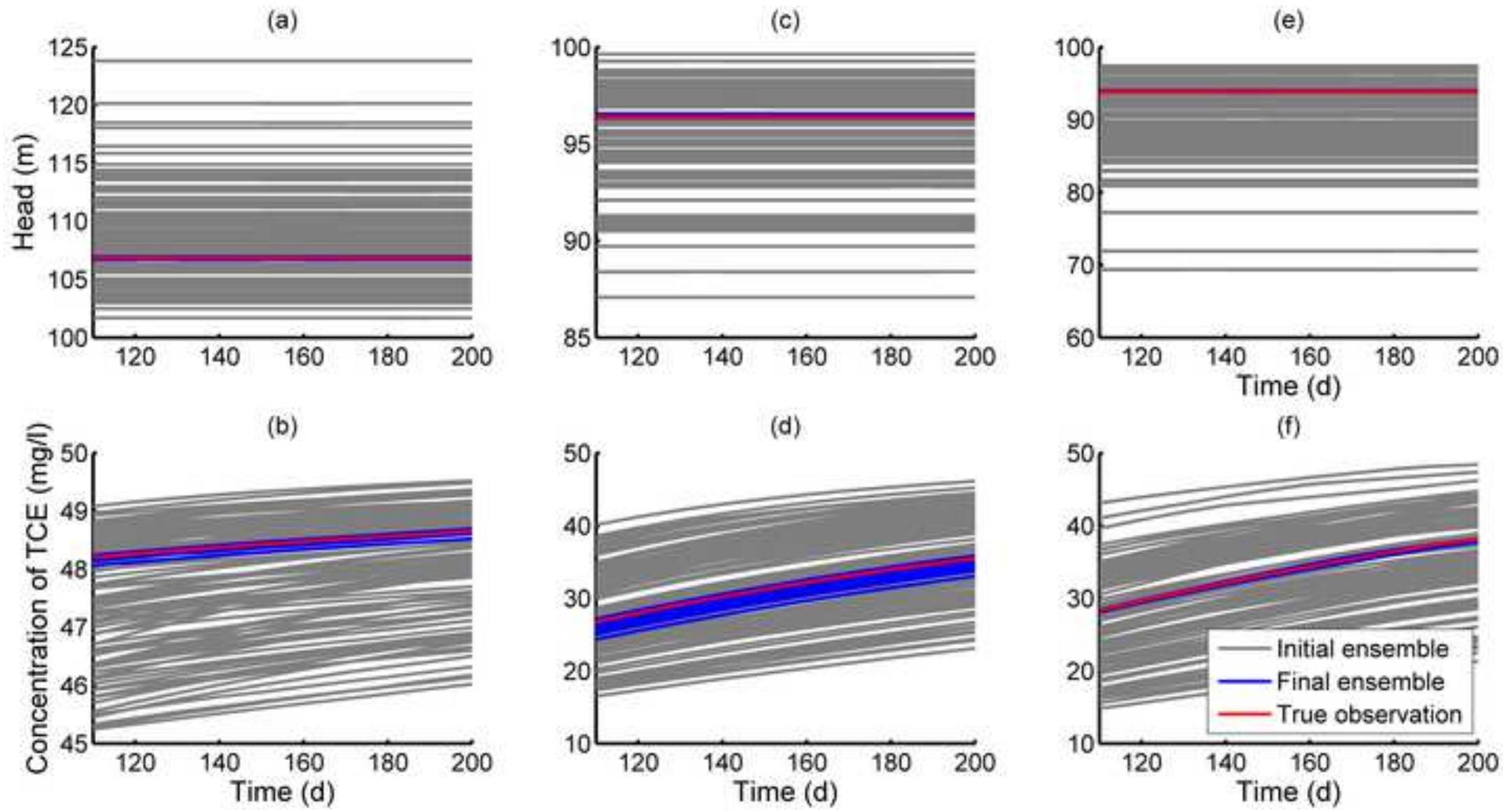












1 Tables

2 **Table 1: Flow and transport parameters used in Case 1**

Flow simulation	Transient state
Total simulation time (days)	10
Stress period	1
Time steps	100
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	150
Model width (m)	5
Model height (m)	5
Starting head (m)	100
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
Injection rate per well (m^3/d)	50
Pumping rate per well (m^3/d)	45

1 **Table 2: Data assimilation related parameters used in different cases**

Case name	Dimension	Number of ensemble (N_e)	Number of assimilation step (N_s)	Optimize or not	Number of optimal sampling locations
Case 1	1	300	10	Y	2
Case 2	2	100	10	Y	2
Case 3	2	50	10	Y	2
Case 4	2	300	10	Y	2
Case 5	2	500	10	Y	2
Case 6	2	1000	10	Y	2
Case 7	2	100	10	Y	1
Case 8	2	100	10	Y	5
Case 9	2	100	10	Y	10
Case 10	2	100	10	Y	20
Case 11	2	100	10	N	(2 fixed)
Case 12	2	100	10	N	(10 fixed)
Case 13	2	100	10	N	(20 fixed)
Case 14	2	100	8(50, 50)*	Y	2
Case 15	2	100	8(50, 50)	Y	5
Case 16	2	100	12(20, 30, 50)	Y	5
Case 17	2	100	12(50, 30, 20)	Y	5

2 * The numbers in the parentheses are the group division of observation time, and the number in front of the
3 parentheses is the number of assimilation steps. For example, 8(50, 50) represents that there are 8 steps in
4 the assimilation and the observation time is divided into two groups with each group having an observation
5 time of 50 days.

1

Table 3: Flow and transport parameters used in Case 2

Flow simulation	Transient state
Total simulation time (days)	100
Stress period	1
Time steps	200
Grid spacing (m)	$5 \times 5 \times 5$
Model length (m)	105
Model width (m)	65
Model height (m)	5
Starting head (m)	50
Porosity	0.3
Specific storage (m^{-1})	0.0001
Longitudinal dispersivity (m)	10
horizontal transverse dispersivity (m)	1
Injection rate per well (m^3/d)	80
Pumping rate per well (m^3/d)	120
TCE injection concentration per well (mg/L)	50

2

3

Table 4: Computational costs

Case name	Times of model invoking	Times of GA invoking	computing time (using the same computer)
Case 2	10	10	2 h
Case 14	8	2	26 min
Case 15	8	2	26 min
Case 16	12	3	25 min
Case 17	12	3	25 min