# Identification of Irregular Non-Point Contaminant Sources Using Ensemble Smoother with Multiple Data Assimilation and Closed Cubic B-Spline Curve Approximation

Wenjun Zhang[a,b,d], Teng Xu[a,b,*], J. Jaime Gómez-Hernández[f], Zi Chen[e], Chunhui Lu[a,c], Guodong Zhang[a,b]

[a]The National Key Laboratory of Water Disaster Prevention, Hohai University, Nanjing, China

[b]College of Water Conservancy and Hydropower Engineering, Hohai University, Nanjing, China

[c]Yangtze Institute for Conservation and Development, Hohai University, Nanjing, China

[d] National Marine Environmental Monitoring Center, 116023 Dalian, China

[e] Nanjing Center, China Geological Survey, Nanjing, China

[f] Institute of Water and Environmental Engineering, Universitat Politècnica de València, Valencia, Spain

* **Corresponding author**.

*E-mail address:* teng.xu@hhu.edu.cn (Teng Xu)

**Abstract:**

Reconstructing the spatial and temporal characteristics of non-point sources is essential for early warning and precise restoration of groundwater pollution. Currently, the identification of non-point contaminant sources is primarily focused on those with simple regular spatial architecture, such as rectangular or circular shapes, while the characterization of irregular non-point contaminant sources has been overlooked, which is attributed to their complex and diverse spatial structures that cannot be readily captured by predefined geometric parameters such as centroid coordinates and side lengths. To address this challenge, we propose a novel approach, an ensemble smoother with multiple data assimilation combined with a closed cubic B-spline curve approximation (ES-MDA-Bs), for the identification of irregular non-point sources. The effectiveness of ES-MDA-Bs in identifying irregular non-point sources in synthetic confined aquifers is assessed and the influence of varying numbers and placements of knot points on its performance for this purpose explored. Then, a comparative analysis between ES-MDA-Bs and the ensemble smoother with multiple data assimilation combined with a rotating ellipse approximation (ES-MDA-RE) ---which has previously demonstrated competency in identifying ellipse-like sources---is performed. The results show that ES-MDA-Bs, with sufficient knot points to approximate an irregular shape, can effectively identify both the irregular spatial structure of the source and its release parameters, i.e., the initial release time, release duration, and mass loading. ES-MDA-Bs clearly outperforms ES-MDA-RE in identifying irregular non-point sources. The accuracy and uncertainty of ES-MDA-Bs are improved as the number of knot points increases.

**Keyword**: Irregular contaminant source identification; Ensemble smoother with multiple data assimilation; Groundwater contamination

## 1. Introduction

Groundwater contamination threatens surrounding ecosystems and may disrupt economic activities. A fundamental approach to mitigate this threat involves tracing the contaminant back to its source and try to accurately reproduce the contaminant plume evolution within an aquifer. Source tracking involves reconstructing its spatial and temporal release characteristics, which are essential for accurate risk assessments and the development of effective remediation strategies (Srivastava and Singh, 2015).

Based on their characteristics, contaminant sources can be classified as either point or non-point. Point sources refer to localized and discrete origins of pollution, often resulting from accidents at specific sites, such as chemical spills or leaks from landfill sites. In contrast, non-point sources are diffuse and continuous, originating from activities that cover large areas, such as fertilizer leaching from agricultural irrigation or effluents from chemical plants. Compared to point sources, non-point sources exert broader geographical and longer-lasting effects on the aquifer, and their spatial distribution often exhibits irregular shapes. Thus, accurately characterizing the irregular spatial structure of a non-point source is a critical factor for effective pollution control.

In the past, due to technical and computational limitations, researchers working on the estimation of source parameters have concentrated in either single point sources (e.g., Hwang and Koerner, 1983; Jha and Datta, 2015; Xu and Gómez-Hernández, 2016, 2018; Chen et al., 2018; Hou et al., 2021; An et al., 2022; Bai and Tahmasebi, 2022; Chang et al., 2022; Ge et al., 2023; Chang et al., 2024) or multiple discrete point sources (e.g., Gorelick et al., 1983; Mahar and Datta, 1997; Hwang et al., 2020; Jamshidi et al., 2020; Anshuman and Eldho, 2022; Kontos et al., 2022; Anshuman and Eldho, 2023; Chen et al., 2023a; Li et al., 2023a; Singh and Mahor, 2023; Wu et al., 2025).

Recently, non-point source identification studies have gradually gained attention. However, the limited studies on non-point source reconstruction have predominantly focused

79   on sources with regular shapes. The existing research on non-point sources can be categorized

80   into two groups based on how the sources are treated during the identification process. The

81   first group involves only identifying the release information parameters of sources, with the

82   spatial structure and location assumed to be known. For instance, in Butera and Tanda (2003)

83   and Pan et al. (2023a), the release history was estimated from a square-shaped source that is

84   uniform in space but changes over time. Both Zhang et al. (2024a) and Zheng et al. (2024)

85   provided accurate estimations of the release strength from a linear contaminant source. Other

86   studies found contaminant release intensities from multiple discrete potential sources with

87   known locations and irregular shapes (e.g., Gzyl et al., 2014; Pan et al., 2021, 2022, 2023b,

88   2025; Li et al., 2023b; Luo et al., 2023). The second group focuses on simultaneously

89   identifying both the release information and the spatial structure of regular-shaped sources.

90   For instance, Mahinthakumar and Sayeed (2005) and Jin et al. (2009) determined the location

91   and shape of a prismatic source by identifying the coordinates of two diagonally-opposed

92   vertices, followed by estimating its release concentration. Mirghani et al. (2009) and Mirghani

93   et al. (2012) reconstructed a cubic areal source by identifying its centroid and side length, and

94   reproduced its release information by calculating the initial concentration. Ayvaz (2016) was

95   the first to try to figure out irregularly shaped sources. He used a new simulation-optimization

96   method to find both the release information and the locations of discrete elements that, when

97   put together, made up the shape of the source. However, when dealing with large-scale non-

98   point sources, this approach may introduce significant errors and computational burdens. Due

99   to the diverse spatial characteristics of non-point sources, inherent uncertainties are

100  substantial, rendering traditional point source identification methods inadequate for

101  addressing non-point source issues. Notably, to the best of our knowledge, current scientific

102  literature lacks effective methodologies capable of systematically characterizing the spatial

103  structure of irregular non-point sources. Hence, a logical and streamlined method is suggested

104 for accurately describing the spatial distribution of irregular sources of any size and shape, as

105 long as some background information about where the contamination came from is known.

106       In our earlier work (Xu et al., 2022), we successfully applied the rotation ellipse approach,

107 which involved approximating a non-point source with an ellipse with five parameters (center

108 point coordinates, semi-major and semi-minor axes, and rotation angle). In this work, we

109 propose a major improvement by shifting to the use of spline curves, which are characterized

110 by a number of knot points, to identify irregular shapes.

111       The ensemble smoother with multiple data assimilation is combined with a closed cubic

112 B-spline curve approximation (ES-MDA-Bs) to deal with irregular non-point source

113 identification problems. It was decided to use the ensemble smoother with multiple data

114 assimilations (ES-MDA) to solve the inverse problem because it has been shown to be good

115 at finding sources in both synthetic (e.g., Xu et al., 2021, 2022) and experimental cases (Chen

116 et al., 2023b). To see how well ES-MDA-Bs can approximate the source's spatial distribution,

117 we try out different configurations by changing the number and placement of knot points. We

118 also perform a comparison between ES-MDA-Bs and the ensemble smoother with multiple

119 data assimilation combined with a rotating ellipse approximation (ES-MDA-RE) to see how

120 well different approximation methods work for finding irregular non-point sources. Note that

121 our study aims to investigate the capability of ES-MDA-Bs to estimate the location, spatial

122 structure, and release information of an irregular non-point source.

123       The structure of this paper is organized as follows: In Section 2, we talk about the

124 algorithms behind both ES-MDA-Bs and ES-MDA-RE in detail. In Section 3, we set up three

125 scenarios to test how well ES-MDA-Bs and ES-MDA-RE work in a confined aquifer. Finally,

126 in Section 4, we compare and analyze the results, and in Section 5, we talk about what needs

127 to be done in the future. Finally, in Section 6, we wrap up the paper with a full summary.

128 **2. Methodology**

129       We propose the ES-MDA-Bs algorithm to identify the parameters of a non-point source

with an irregular spatial structure. To demonstrate the effectiveness of ES-MDA-Bs, we use

the ES-MDA-RE algorithm for comparison. The computational procedures of both algorithms

are described in the following sections.

**2.1 Ensemble smoother with multiple data assimilation combined with a closed cubic B-spline curve approximation (ES-MDA-Bs)**

ES-MDA-Bs is a data assimilation method designed to identify complex source with a

spatial irregular shape. It integrates two key components: ES-MDA for parameter

identification and a closed B-spline curve to delineate the source shape through a number of

knot points. This approach assimilates concentration observations to determine the

coordinates of the knot points on the spline curve and the source release parameters, enabling

the reconstruction of both the spatial structure and release history of the source.

**2.1.1 Ensemble smoother with multiple data assimilations (ES-MDA)**

ES-MDA is an ensemble-based data assimilation approach proposed by Emerick and

Reynolds (2013), that incorporates an iterative scheme to handle the non-linear state equations

in hydrological modeling. It is an evolution of the ensemble smoother (Evensen and Van

Leeuwen, 2000), which, in turn, evolves from the ensemble Kalman filter (Li et al., 2012).

This approach assimilates observations of state variables from all time steps into the ensemble

of model simulations at once. The computational process can be typically divided into two

key steps: forecast and update.

During the forecast step, a forward model, represented by equation $\Psi(\bullet)$, forecasts the

state variables $C_j^f$ (in our case, solute concentrations) for all time steps at the $j^{th}$

assimilation iteration, using both the initial state variables $C_0$ and the model parameters (in

our case, the parameters defining the source location and release history) determined during

the previous iteration $P_{j-1}^a$, as model inputs. Note that forecasting from $C_0$ implies that each

forward model simulation is run from the same initial state at time 0.

$$C_j^f = \Psi(C_0, P_{j-1}^a) \tag{1}$$

156    During the update step, observation data from all time steps are assimilated to refine the

157    model parameters. The assimilation process adjusts the model parameters using the

158    discrepancy between forecasted states at observation locations $C_j^{f,o}$ and observation data in

159    the same wells. The objective is to obtain an updated parameter vector $P_j^a$ at the $j^{th}$

160    iteration based on the parameter state from the previous iteration $P_{j-1}^a$ aimed to reducing the

161    discrepancies between forecasts and observations,

162
$$P_j^a = P_{j-1}^a + K_j(C^o + \sqrt{a_j}\varepsilon_j - C_j^{f,o}) \tag{2}$$
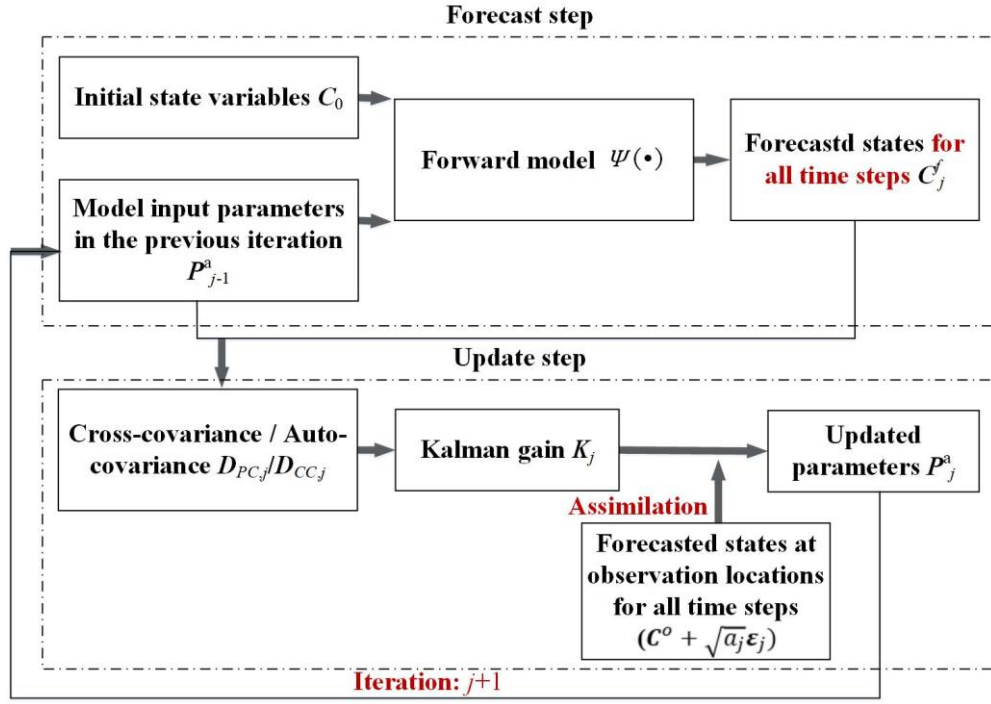
163    with

164
$$K_j = D_{PC,j}(D_{CC,j} + a_j R)^{-1} \tag{3}$$

165    where $C^o + \sqrt{a_j}\varepsilon_j$ is the observed data already affected by an observation error $\varepsilon_j$, which is

166    damped by factor $a_j$. According to Emerick and Reynolds (2013), the damping factors should

167    be chosen as to satisfy $\sum_{j=1}^{Na}\frac{1}{a_j} = 1$ and in decreasing order following an arithmetic or a

168    geometric regression (specifically, $a_j$ is equal to the total number of assimilation iterations

169    $Na$ for all iterations). $K_j$ is the Kalman gain and depends on the error covariance $R$, the

170    cross-covariance $D_{PC,j}$ between parameters and states at the observation locations at all time

171    steps and the auto-covariance $D_{CC,j}$ between states at observation locations at all time steps

172    at the $j^{th}$ iteration.

173        In this study, we evaluate the performances of ES-MDA for two different

174    implementations: The ellipse-base method and the spline-base one. The source parameters to

175    be identified $P$ will be the initial release time $T^0$ [T], the release duration $\Delta T$ [T], the

176    constant release mass-loading $M$ [MT$^{-1}$], and the geometric parameters $Z$ ($Z = Z_s$ or $Z_e$, Z

177    varies depending on the source shape approximation method, as will be explained in the

178    following subsections).

179
$$P = [Z, T^0, \Delta T, M]^T \tag{4}$$

180    ES-MDA serves as the core inversion algorithm in this study. The workflow is

181    summarized in Figure 1, where the source parameters $P$ vary according to the source shape

182    approximation method employed.



183

184    **Figure 1**: Workflow of ES-MDA.

185    **2.1.2 Closed cubic B-spline curve approximation**

186    The B-spline curve, initially proposed by Isaac Jacob (Schoenberg, 1946), has gained

187    significant popularity in computer-aided design and numerical analysis due to its versatility

188    (e.g., De Boor and De Boor, 1978; Piegl and Tiller, 1996; Park and Lee, 2007; Farin, 2014).

189    The curve, characterized by degree $k$, is expressed as the sum of basis functions $B_{ia,k}(t)$

190    multiplied by control points $A_{ia}$. The curve within data point segment $[D_i, D_{i+1})$ can be

191    written as

$$Z_{S_i}(t) = \sum_{ia=i}^{i+k} B_{ia,k}(t) A_{ia} \tag{5}$$

193    where $ia$ is the minimum number of control points required to ensure the smoothness of the

194    curve; $i$ is the number of knot points, $i = 0, \ldots, n-1$. The basis function $B_{ia,k}(t)$ are

195    given by

$$B_{ia,0}(t) = \begin{cases} 1 & t \in [D_i, D_{i+1}) \\ 0 & t \in otherwise \end{cases} , k = 0 \tag{6}$$

197     and

$$B_{ia,k}(t) = \frac{1}{k!}\sum_{j=0}^{i+k-ia}(-1)^j \binom{k+1}{j}(t + i + k - ia - j)^k, k > 0 \qquad (7)$$

199     where $\binom{k+1}{j}$ is the binomial coefficient representing the number of combinations of $j$

200     elements that can be chosen from a set of $(k + 1)$ distinct elements. The cubic B-spline

201     curve $(k = 3)$ offers a good balance between smoothness and computational efficiency,

202     making it well-suited for a wide range of applications.

203         By substituting $k = 3$ into Eq. (5), the cubic B-spline is given by

204     $Z_{S_i}(t) = \frac{1}{6}(1 - t)^3 A_i + \frac{1}{6}(3t^3 - 6t^2 + 4)A_{i+1} + \frac{1}{6}(-3t^3 + 3t^2 + 3t + 1)A_{i+2} + \frac{1}{6}t^3 A_{i+3}$

205                                                                                  (8)

206     where the coordinates of $A_i$, $A_{i+1}$, $A_{i+2}$ and $A_{i+3}$ are computed to ensure that the curve

207     passes through the knot points $D_i$ and $D_{i+1}$, and that there is first and second derivative

208     continuity at these points. For readers unfamiliar with splines and the distinction between

209     control points and knot points: To elaborate, control points are used to construct the curve

210     and establish the relationship between the knot points and the resulting B-spline curve (see

211     Eq. (5)). Conversely, knot points represent specific locations that the curve must pass through.

212     The relationship between control points and knot points can be formally described in Eq. (9).

213     For further details on fitting knot points to a curve, please refer to the following references

214     (e.g., de Boor, 1978; Lyche and Mørken, 1999; Juhász and Hoffmann, 2004).

215     $$D_i = \frac{1}{6}(A_i + 4A_{i+1} + A_{i+2}) \qquad (9)$$

| Algorithm 1: ES-MDA-Bs | |
| --- | --- |
| 1: | Generate initial ensembles of source parameters ($P = [Z_s, T^0, \Delta T, M]^T$). |
| 2: | for $j = 1:N_a$ |
| 3: | Forecast concentrations in the observation wells ($C_j^{f,o}$). |
| 4: | Update source parameters ($P$) based on the misfits between forecasts and observations ($C^o + \sqrt{a_j}\varepsilon_j - C_j^{f,o}$). |
| 5: | end for |

216

217    Consequently, the core steps of the ES-MDA-Bs algorithm are summarized in Algorithm

218    1. Note that the spatial shape of the source location ($Z_s$) is defined by a set of $n$ knot points

219    to be identified

220

$$Z_s = [\{D_i\}, i = 0, \ldots, n-1]^T \tag{10}$$

221 **2.2 Ensemble smoother with multiple data assimilation combined with a rotating ellipse**
222 **approximation (ES-MDA-RE)**

223    In our previous work (Xu et al., 2022), we proposed ES-MDA-RE to identify irregular

224    shapes that could be approximated by an ellipse. The assimilation process is similar to ES-

225    MDA-Bs, with the primary difference that ES-MDA-RE uses a different parameterization of

226    the contaminant source shape. While in ES-MDA-Bs the unknowns are the coordinates of the

227    $n$ knot points, in ES-MDA-RE the unknowns ($Z_e$) are the parameters defining the ellipse: the

228    center point coordinates $(Xs, Ys)$ [L, L], the lengths of the semi-major and semi-minor axes

229    $Ra$ [L] and $Rb$ [L], and the clockwise rotation angle with respect to the $x$-axis $B$ [°].

230

$$Z_e = [Xs, Ys, Ra, Rb, B]^T \tag{11}$$

| Algorithm 2: ES-MDA-RE |
| --- |
| 1:    Generate initial ensembles of source parameters ($P = [Z_e, T^0, \Delta T, M]^T$). |
| 2:    for $j = 1:N_a$ |
| 3:        Forecast concentrations in the observation wells ($C_j^{f,o}$). |
| 4:        Update source parameters ($P$) based on the mists between forecasts and observations $(C^o + \sqrt{a_j}\varepsilon_j - C_j^{f,o})$. |
| 5:    end for |

231

232 **3. Application**

233    We construct a two-dimensional synthetic confined aquifer, discretized into a grid of

234    80×80×1 cells, with each cell size measuring 10 [L]×10 [L]×80 [L]. (No specific units will

235    be used; any set of consistent units will result in the results presented.) In this study, the release

236    of a non-reactive contaminant is modeled under a steady-state groundwater flow condition,

237    and solute movement is driven solely by advection and dispersion. Consequently, the

238  governing equations for the forward model include the steady-state groundwater flow based

239  on Darcy's law and the principle of continuity,

240
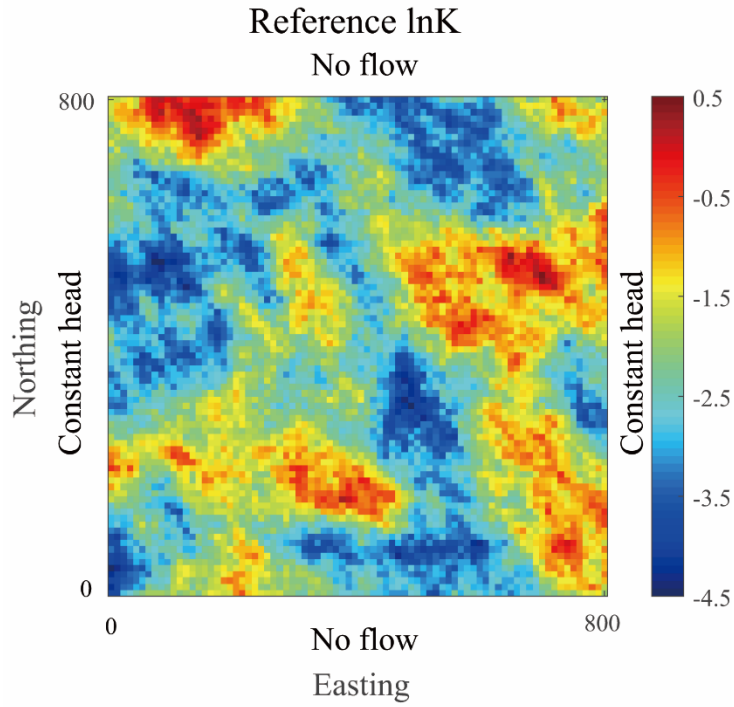$$Q = -\nabla \cdot (K\nabla H) \qquad (12)$$

241  where $Q$ represents sources and sinks per unit volume $[T^{-1}]$; $\nabla \cdot$ is the divergence operator;

242  $\nabla$ is the gradient operator; $K$ is the hydraulic conductivity $[LT^{-1}]$; and $H$ is the hydraulic

243  head $[L]$; and the transport equation,

244
$$\frac{\partial(\theta_e C)}{\partial t} = \nabla \cdot [\theta_e(\alpha_m + \beta v) \cdot \nabla C] - \nabla \cdot (\theta_e v C) - q_c C_s \qquad (13)$$

245  where $C$ is the contaminant source concentration $[ML^{-3}]$; $t$ is the simulation time $[T]$; $\alpha_m$

246  is the molecular diffusion coefficient $[L^2 T^{-1}]$; $\beta$ is the dispersivity tensor $[L]$; $q_c$ is the

247  volumetric flow rate per unit cross-section $[LT^{-1}]$; and $C_s$ is the concentration of the sources

248  or sinks $[ML^{-3}]$; $\theta_e$ [-] is the effective porosity and $v$ $[LT^{-1}]$ is the velocity given by $v =$

249  $(-K\nabla H)/\theta_e$.

250      To simulate groundwater flow and solute transport processes, the forward model requires

251  specific input parameters to solve for the state variables. Figure 2 presents the natural

252  logarithm of the hydraulic conductivity field, which exhibits a heterogeneous Gaussian

253  distribution generated using the GCOSIM3D program, a multivariate multi-Gaussian

254  sequential simulation code (Gómez-Hernández and Journel, 1993). The Gaussian distribution

255  is defined by the parameters listed in Table 1. The boundary conditions for the aquifer are as

256  follows: the east and west boundaries are set to constant heads of 80 [L] and 300 [L],

257  respectively, while the north and south boundaries are impermeable. Porosity is homogeneous

258  and equal to 0.3 [-]. For the solute transport process, the initial concentration across the

259  domain is set to zero, and the dispersion coefficient $\beta$ is assumed to be constant in both space

260  and time. The total transport simulation time is 10950 [T], divided into 100 stress periods of

261  109.5 [T] each. The non-reactive contaminant source releases at a constant mass-loading rate

262  from the 10th time step (approximately 985.5 [T]) to the 30th time step (around 3285.0 [T]),

263  resulting in a release duration of 21-time steps (approximately 2299.5 [T]), as detailed in

264    Table 2. The observational errors are assumed to follow a normal distribution with zero mean

265    and a variance of 0.01. Groundwater flow and solute transport processes are simulated using

266    the MODFLOW and MT3D programs, respectively.



267

268    **Figure 2**: Reference log-conductivity field. An indication of the type of boundary conditions

269    used for the solution of the flow equation is also shown.

270

271    **Table 1.** Parameters used to generate the lnK field

|      | Mean | Std.dev. | Variogram | $\lambda_{max}$ | $\lambda_{min}$ | Angle |
|------|------|----------|-----------|-----------------|-----------------|-------|
| lnK  | -2   | 1        | Spherical | 300             | 200             | 135   |

272

273                    **Table 2.** Source release parameters

| Parameters | Values |
|------------|--------|
| Initial release time $[T^0]$ | 985.5 |
| Release duration time $[\Delta T]$ | 2299.5 |
| Mass loading rate $[M]$ | 100 |

274        As mentioned, this study focuses on identifying non-point source with complex irregular

275    spatial distribution. The reference shape for the source is presented in Figure 3 (a). To

–12–

276 demonstrate the performance of ES-MDA-Bs in identifying irregular non-point sources, three

277 scenarios are developed: in scenario S1 and S2, the spatial structure of the source is

278 approximated using 10 and 5 knot points, respectively, which form a continuous path based

279 on the initial position vector $D_i$ (Eq.10) and are identified by ES-MDA-Bs; and in scenario

280 S3, the source is assimilated to an ellipse and its 5 defining parameters identified using ES-

281 MDA-RE. Scenarios S1 and S2 are designed to evaluate the impact of the number of knot

282 points on the approximation of the irregular source spatial structure using ES-MDA-Bs.

283 Scenario S3 serves to compare the performance of the ellipse-based method and the spline-

284 based one, eventually highlighting the advantage of using ES-MDA-Bs over ES-MDA-RE

285 for irregular source identification. In all scenarios, the number of assimilation steps is also

286 tested. The basic parameters defining the three scenarios are listed in Table 3.



(a)    (b)

288 **Figure 3**: Reference contamination source (a) and the computed envelope of generated shapes

289 within the suspect range for knot points and elliptical control parameters, used to initialize the

290 parameters prior to the start of the assimilation process (b). The red triangles are the

291 observation locations. The black squares are the validation locations (not used during the

292 assimilation process). The suspects are in (b) correspond to scenarios S1 (black), S2 (green)

293 and S3 (red).

294 For each scenario, we generate an ensemble of 500 realizations, consisting of 23-tuples

295 for S1, 13-tuples for S2, and 8-tuples for S3. Within each tuple, the values are randomly drawn

296 from a uniform distribution and used as initial inputs for the inversion process. The parameter

297 ranges are broad to account for significant uncertainty and are chosen from uniform

298 distributions as indicated next: $x$-coordinates of knot points $Xs \in U[80, 380]$, y-coordinates

299     of knot points $Ys \in U[400, 700]$, initial release time $T^0 \in U[0, 3175.5]$, release duration

300     $\Delta T \in U[1204.5, 6679.5]$ and mass-loading rate $M \in U[90, 150]$, $x$-coordinate of ellipse

301     center $Xsc \in U[160, 200]$, $y$-coordinate of ellipse center $Ys \in U[480, 580]$, semi-major

302     axis $Ra \in U[110, 210]$, semi-minor axis $Rb \in U[50, 120]$ and clockwise rotation angle

303     $B \in U[0, 90]$. The prior suspected parameter ranges are summarized in Table 3. The envelops

304     of all 500 realizations of initial guesses of the contaminant source are shown in Figure 3 (b).

305     It is important to recall that the update step is non-convex, meaning that the updated parameter

306     values can be outside the initial uncertainty intervals if the dynamics of the system indicate

307     that the source was initially clearly off the reference.

308                              **Table 3.** Definition of scenarios

| Scenario | | S1 | | | | S2 | | | | S3 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Number of assimilation iterations [$la$] | | 0 | 1 | 4 | 7 | 0 | 1 | 4 | 7 | 0 | 1 | 4 | 7 |
| Approximation | | Spline curve | | | | | | | | Rotating ellipse | | | |
| Number of knots [$lk$] | | 10 | | | | 5 | | | | / | | | |
| Suspect Range | $x$-coordinate of knot [$Xs$] | 80-380 | | | | | | | | / | | | |
| | $y$-coordinate of knot [$Ys$] | 400-700 | | | | | | | | / | | | |
| | $x$-coordinate of center point of ellipse [$Xsc$] | / | | | | | | | | 160-200 | | | |
| | $y$-coordinate of center point of ellipse [$Ysc$] | / | | | | | | | | 480-580 | | | |
| | Semi-major axis of ellipse [$Ra$] | / | | | | | | | | 110-210 | | | |
| | Semi-minor axis of ellipse [$Rb$] | / | | | | | | | | 50-120 | | | |
| | Clockwise rotation angle [$B$] | / | | | | | | | | 0-90 | | | |
| | Initial release time [$T^0$] | 0-3175.5 | | | | | | | | | | | |
| | Release duration time [$\Delta T$] | 1204.5-6679.5 | | | | | | | | | | | |
| | Mass-loading rate [$M$] | 90-150 | | | | | | | | | | | |

309

## 4. Results

Figure 4 shows the probability of the source location in all three scenarios before any updates and after $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ assimilation iterations. The statistic $Pr_i$ is employed to quantify the probability of the source occurring in any given unit $i$ within the aquifer, calculated as the ensemble mean of an indicator $Ir_{j,i}$.

$$Pr_i = \frac{1}{N}\sum_{j=1}^{N} Ir_{j,i} \tag{14}$$

where $N$ denotes the number of realizations included in the ensemble; the indicator $Ir_{j,i}$ denotes the presence of the source in cell $i$ for the $j^{th}$ realization, where a value of 1 indicates the source is present in cell $i$, and 0 indicates its absence. The analysis of Figure 4 reveals a significant level of uncertainty in the initial source distribution for all scenarios, which progressively diminishes with each assimilation iteration, reaching a minimum in the final step. Ultimately, S1 and S2 show clear irregular spatial structures of the potential source, while S3 does not. These findings demonstrate that the spline curve is more effective in approximating the spatial distributions of irregular non-point sources than the rotating ellipse, though its accuracy may be limited by the number of knot points.

**Figure 4**: Probability of source location as computed from the source positional parameters updated at the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ assimilation iterations in scenarios S1-S3. The left column corresponds to ES-MDA-Bs with 10 knot points, the center column corresponds to ES-MDA-Bs with 5 knot points, the right column corresponds to ES-MDA-RE.
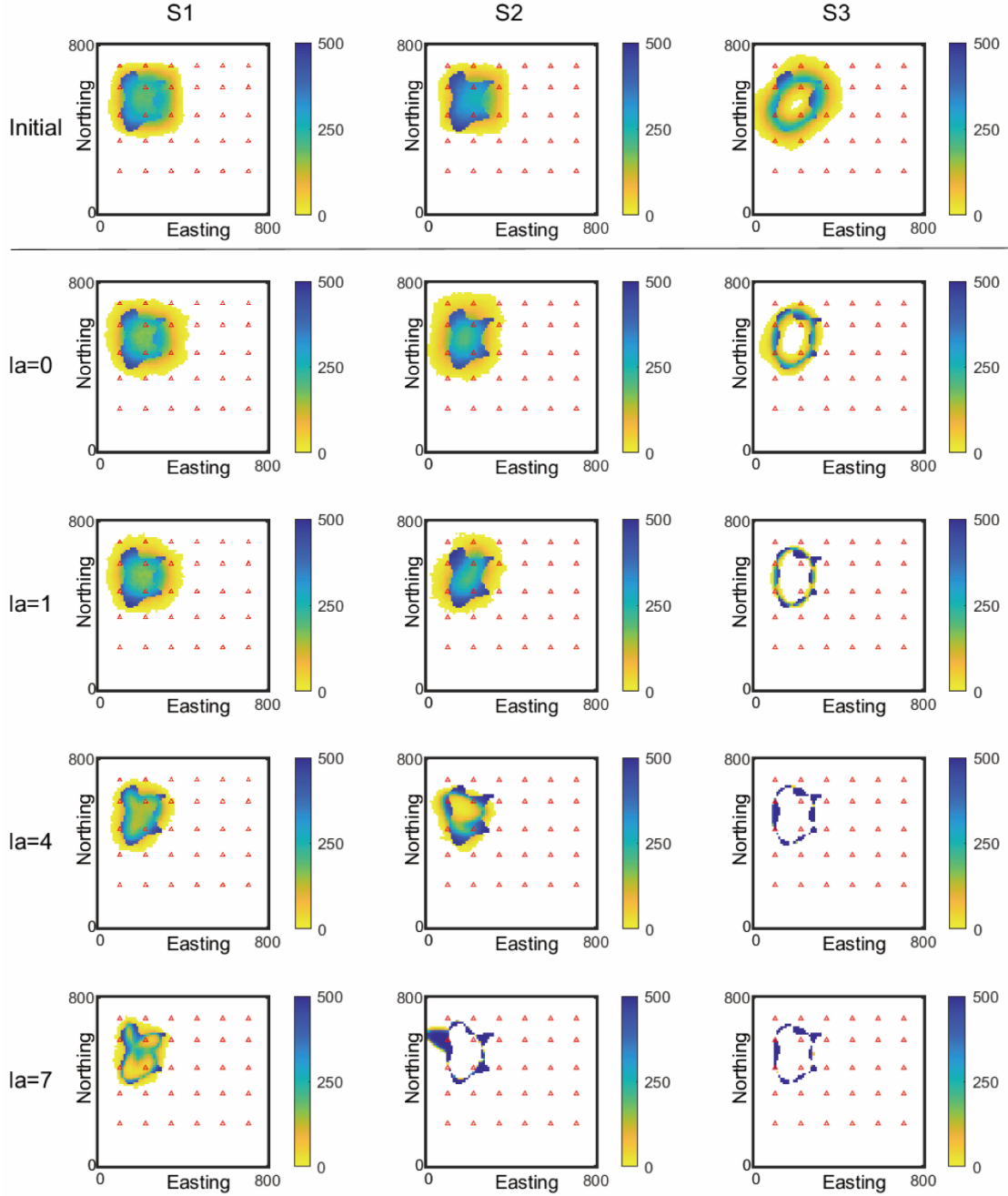
Figures 5 presents a statistic assessing the accuracy of location and shape updates of the potential source through assimilating concentration observations at the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ iterations across all three scenarios. This statistical measure, the Total Ensemble Absolute Bias ($TEAB$), quantifies accuracy across the domain by integrating the absolute deviations of the estimated source indicator $Ir_{j,i}$, from the reference indicator $Ir_{ref,i}$, for each spatial unit

335    $i$. Specifically, $TEAB$ for any unit $i$, $TEAB_i$, is computed as:

336                         $$TEAB_i = \sum_{j=1}^{N} |Ir_{j,i} - Ir_{ref,i}| \qquad (15)$$

337    where $Ir_{ref,i}$ denotes the presence of reference source in cell $i$. Note that an integral value

338    of 0 indicates a fully accurate reproduction of source location and spatial structure, while a

339    value of 500 means that the source does not appear in any realization. The analysis of Figure

340    5 shows a gradual convergence of the potential source structure toward the reference spatial

341    distribution with increasing assimilation iterations in S1. After the 7[th] iteration, minor

342    deviations appear at high-curvature edges in the updated source locations. These deviations

343    cause the updated positions to shift inward, leading to a slight underestimation of the spatial

344    distribution extent of the source. In contrast, S2 shows minimal reduction in deviation

345    between estimated and true indicator values near the source boundary over successive

346    iterations, with a noticable northwest misfit emerging on areas of marked curvature, resulting

347    in an inaccurate restructuring of the source spatial distribution. These findings suggest that

348    when using ES-MDA-Bs with limited-configured knot points for source localization,

349    significant misjudgments occur in high-curvature areas, despite its capacity to estimate

350    irregular spatial structures. Incorporating additional knot points enables ES-MDA-Bs to

351    perform optimally, though it tends to slightly underestimate the spatial extent. Unlike S1 and

352    S2, S3 captures source locations in elliptical areas, with restricted resolution of the

353    irregularities.

**Figure 5**: $TEAB$ computed with the initial and updated ensembles of positional parameters after the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ data assimilation iterations in scenarios S1-S3.
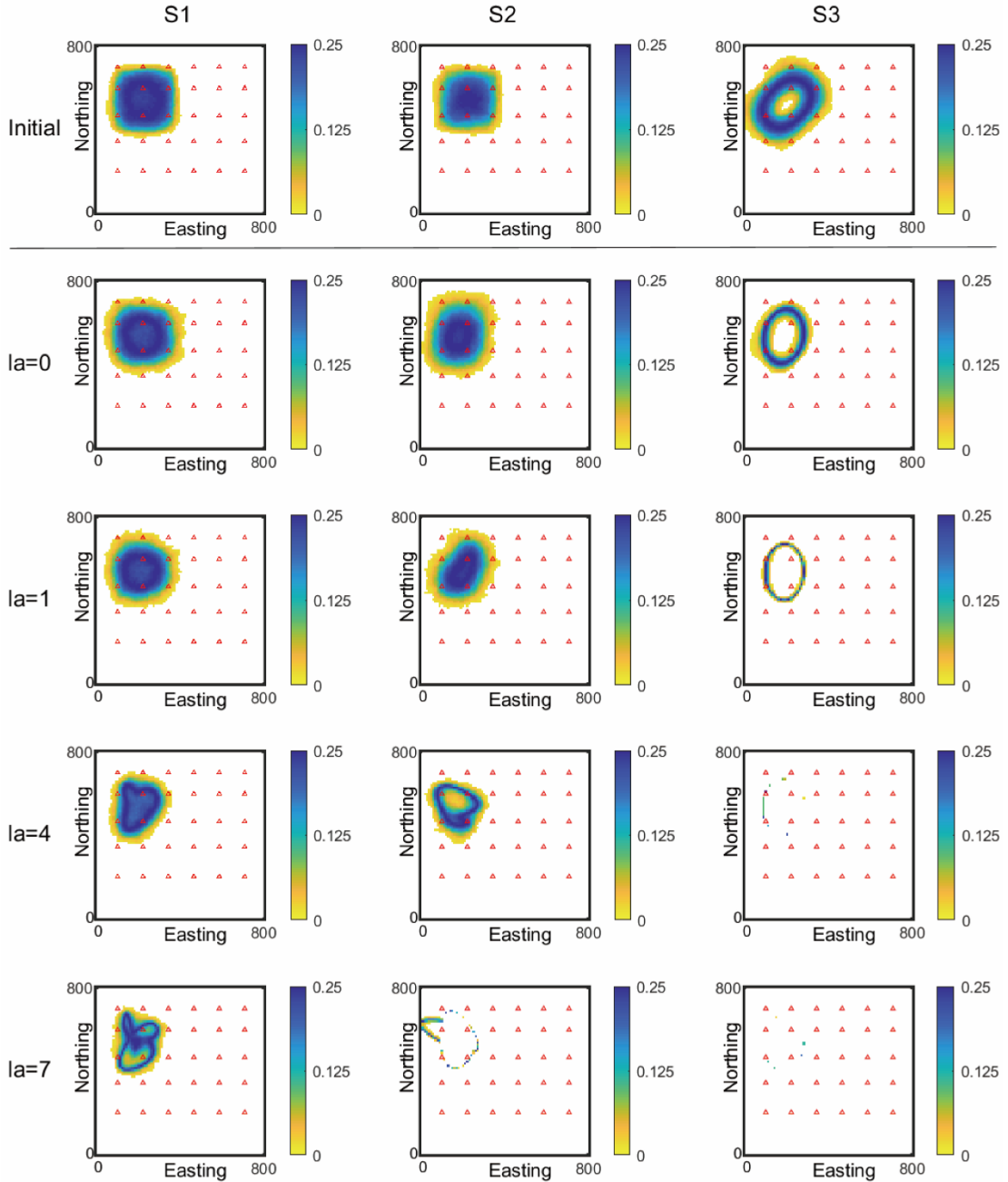
Figure 6 presents the ensemble variance of source indicator values for all three scenarios, providing a measure of variability in the indicator within the ensemble at each potential source location $i$, that can be written as:

$$EnsVar_i = \sigma_{Ir_i}^2 \tag{16}$$

where $\sigma_{Ir_i}$ denotes the ensemble standard deviation of source indicator at location $i$.

Figure 6 shows a significant reduction in initial spatial variability in all three scenarios as the

363    number of iterations increases, reaching a minimum after the 7th iteration. In the final step, S1

364    displays the largest variance along the edges of the source, which is reasonable. Conversely,

365    S2 and S3 show a minimal spatial variability but lacked accuracy, suggesting a severe

366    underestimation of uncertainty. These findings demonstrate that, when using ES-MDA-Bs to

367    estimate the spatial structure of an irregular source, increasing the number of identified knot

368    points can mitigate uncertainty underestimation and improve spatial accuracy.



369

370    **Figure 6**: Ensemble variance computed with the initial and updated ensembles of positional

371    parameters after the 0th, 1st, 4th, and 7th data assimilation iterations in scenarios S1-S3.

372

373      To further clarify the influence of the number of knot points on the reconstruction

374      accuracy and ensemble variability in identifying source location and spatial structure, we

375      evaluated the sensitivity of domain-integrated $TEAB$ and ensemble variance to the number of

376      spline knots. Figure 7 displays the domain-integrated $TEAB$ and ensemble variance computed

377      using ES-MDA-Bs with 5, 8, 10, and 12 knots after the 7th iteration. As the number of knots

378      increases, uncertainty initially rises from a minimum level, while reconstruction accuracy

379      improves and stabilizes around 10 knots but deteriorates sharply beyond this point. This trend

380      is attributed to excessive spline flexibility, which increases sensitivity to local oscillations

381      caused by ensemble uncertainty in knot placement during data assimilation, potentially

382      inducing spatial distortion. Although the spatial extents of source structures reconstructed

383      with 8 and 10 knots are similar, ensemble variance grows with knot count, reflecting an

384      enhanced capacity of the algorithm to explore a wider solution space. These results indicate

385      that insufficient knot configurations restrict the ability of ES-MDA-Bs to capture structural

386      variability, leading to underestimation of uncertainty. However, increasing knot density to

387      mitigate this limitation introduces a critical trade-off between uncertainty representation and

388      reconstruction fidelity. Therefore, in this study, an ES-MDA-Bs configuration with 10 knots

389      achieves reliable performance in estimating both the location and spatial structure of irregular

390      sources. To develop general guidelines for optimal knot configuration, future work will

391      investigate adaptive strategies for determining both knot number and placement.



392

**Figure 7**: Sensitivity analysis of the integrated $TEAB$ and ensemble variance with respect to the number of spline knots, based on the updated ensembles of positional parameters obtained after the 7th data assimilation iteration using ES-MDA-Bs.

As mentioned above, when limited-configured knot points are used for approximation, a clear misfit with underestimated uncertainty appears due to filter inbreeding. This issue stems from ES-MDA-Bs, as an ensemble-based data assimilation algorithm, suffering from spurious correlations in ensemble covariance when the ensemble size is smaller than the number of assimilated observations. To address it, we employed localization technique to the ensemble-based algorithm, with detailed explanations provided in our previous work (Zhang et al., 2024b). Here, we evaluate the capability of ES-MDA-Bs combined with localization (LES-MDA-Bs) in S2 to improve the accuracy of spatial structure identification. Figure 8 presents three statistics, probability, $TEAB$ and ensemble variance of the location and shape of the potential source in S2 using LES-MDA-Bs before and after 0th, 1st, 4th, and 7th assimilation iterations. As the assimilation iterations increase, the potential source gradually aligns with the reference spatial extent, with obvious deviations and uncertainty underestimation removed. This result indicates that while LES-MDA-Bs with limited-configured knot points can capture the general extent of an irregular source, though it struggles to reconstruct the detailed spatial structure.

**Figure 8**: Probability, $TEAB$ and ensemble variance computed with the initial and updated ensembles of positional parameters after the 0[th], 1[st], 4[th], and 7[th] data assimilation iterations in scenario S2 using LES-MDA-Bs.
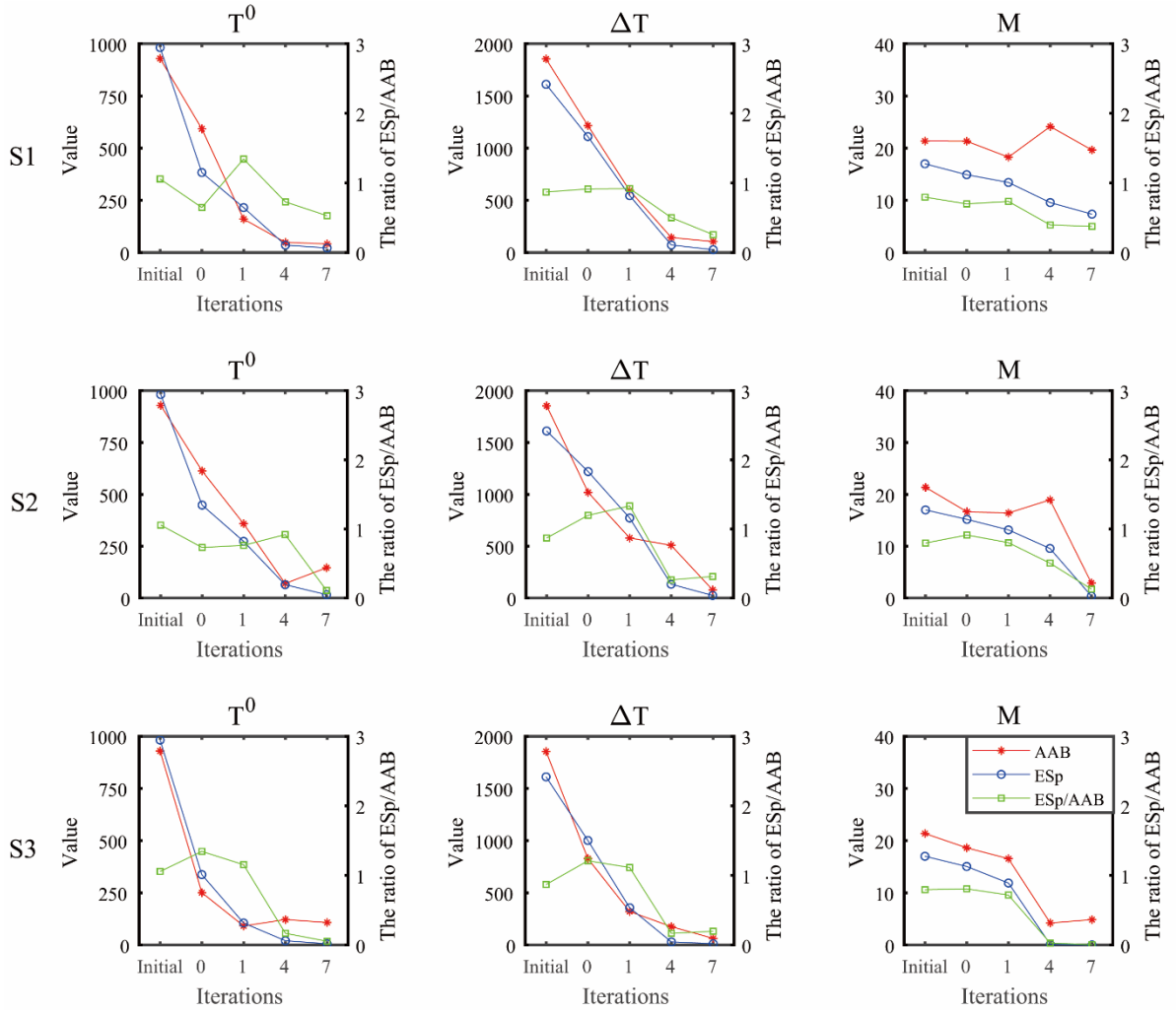
Figures 9 and 10 show three statistical measures, the average absolute bias ($AAB$), the ensemble spread ($ESp$), and boxplot, of the release information parameters ($T^0$, $\Delta T$, $M$), in all three scenarios before and after 0[th], 1[st], 4[th], and 7[th] updating iterations. Specifically, the $AAB$, which quantifies the accuracy of the release parameters, can be written as:

$$AAB = \frac{1}{N}\sum_{j=1}^{N}|Pr_j - Pr_{ref}| \tag{17}$$

420      where $Pr_j$ denotes the source release information value for the $j^{th}$ realization; $Pr_{ref}$ is the

421      corresponding reference value. And the $ESp$, measuring the variability within the ensemble
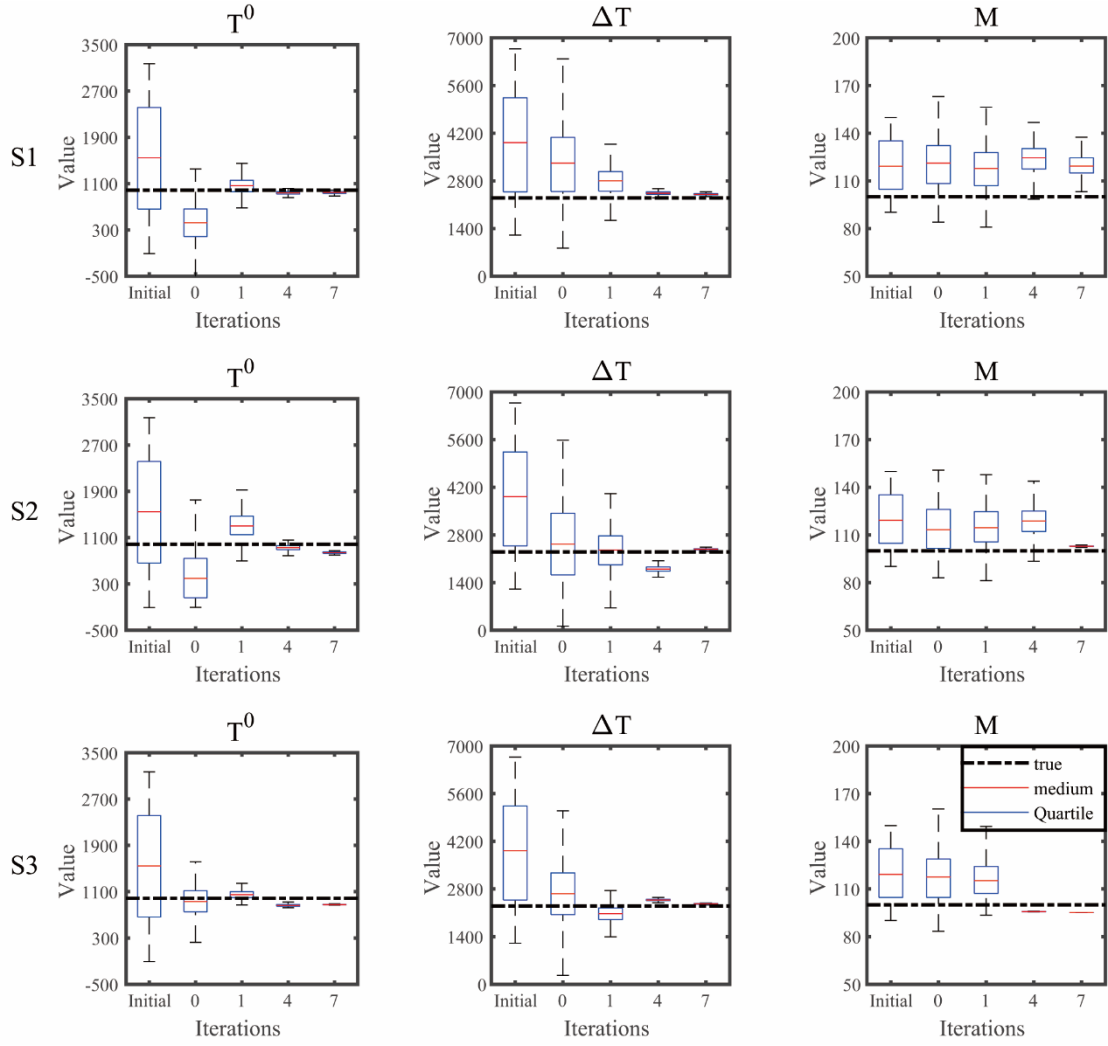
422      for release parameters, can be written as:

423
$$ESp = \sqrt{\sigma_{Pr}^2} \tag{18}$$

424      where $\sigma_{Pr}^2$ denotes the ensemble variance of the release information value. When the ratio

425      of $ESp/AAB$ equals to 1, indicates a satisfactory performance of the smoother (Xu et al.,

426      2013, 2022). Figure 9 shows a progressive decline in $AAB$, $ESp$ and the ratio of $ESp/AAB$

427      with each assimilation, reaching their minimum values after the 7th iteration. Specifically, for

428      scenario S1, the *AAB* of the mass-loading rate exhibits a slight increase after the 4th iteration

429      but subsequently decreases to a lower level by the final update. This temporary fluctuation

430      may result from localized ensemble perturbations introduced during the assimilation process,

431      which can occasionally lead to variations in intermediate estimates. However, the overall

432      trend remains stable, as indicated by the *ESp/AAB* ratio of the mass-loading rate, which

433      remains close to 1 throughout the iterations. This suggests that the algorithm consistently

434      performs well in estimating the mass-loading rate, despite minor fluctuations during

435      individual updates. Additionally, it can be observed that the *ESp/AAB* ratio in S1 stabilizes

436      closer to 1 after the final update compared to that in S3, indicating that ES-MDA-Bs

437      demonstrates greater stability in identifying source information. However, a slight bias is

438      observed in the final release mass-loading rate of S1 compared to the reference value. To

439      further clarify this bias, the boxplots have been shown in Figure 10. This figure illustrates

440      how ensemble medians of the release parameters progressively align with reference values,

441      with uncertainty decreasing to a minimum in the final step. However, due to the

442      underestimation of source spatial extent as mentioned, the final mass-loading rate in S1

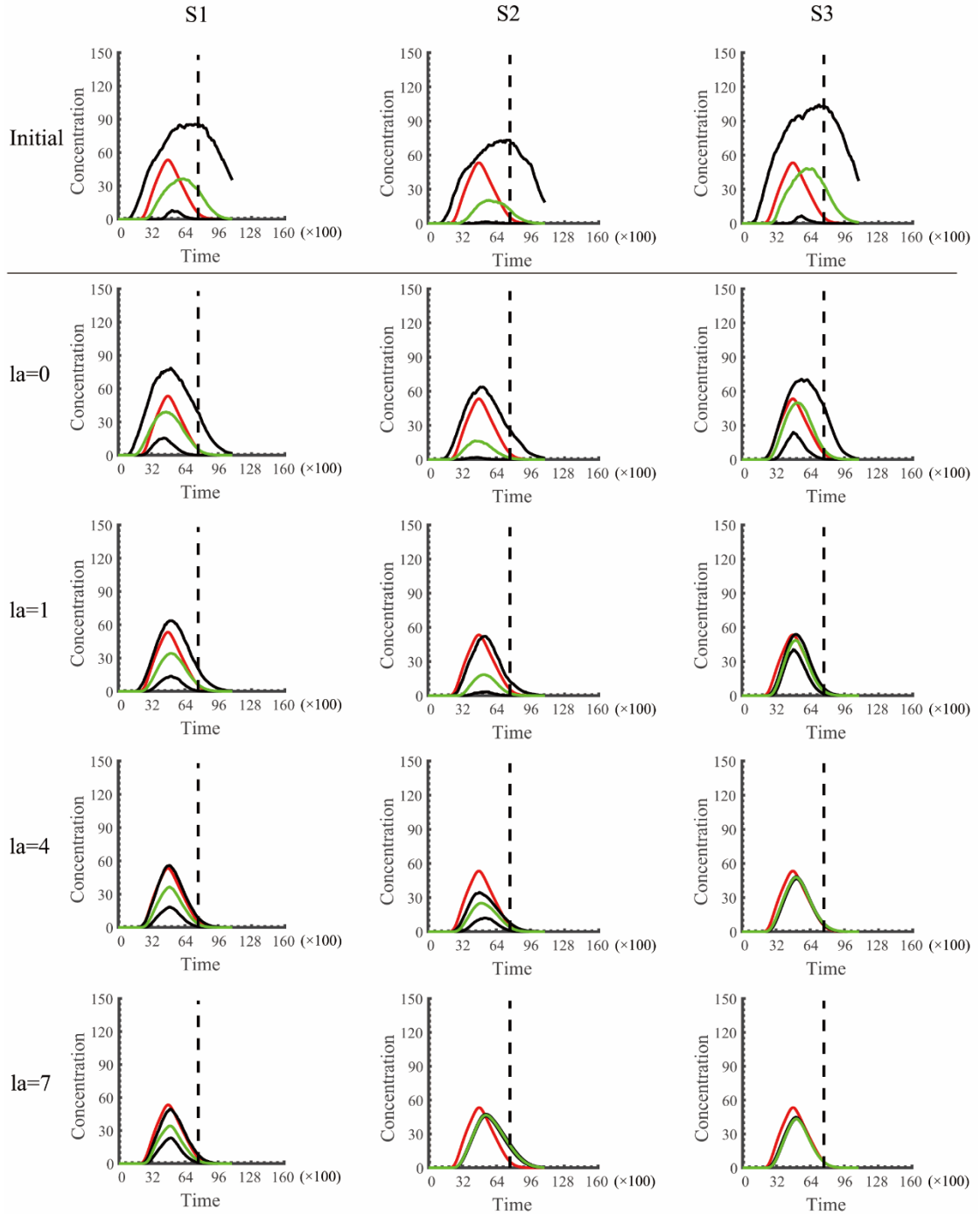443      remains overestimated, consistent with the mass conservation principle.

**Figure 9**: $AAB$, $ESp$ and the ratio of $ESp/AAB$ computed with the initial and updated ensembles of release information parameters of contaminant source, including $T^0$, $\Delta T$ and $M$, after the 0th, 1st, 4th, and 7th data assimilation iterations in scenarios S1-S3. The red line corresponds to $AAB$, the blue line corresponds to $ESp$ and the green line corresponds to $ESp/AAB$.
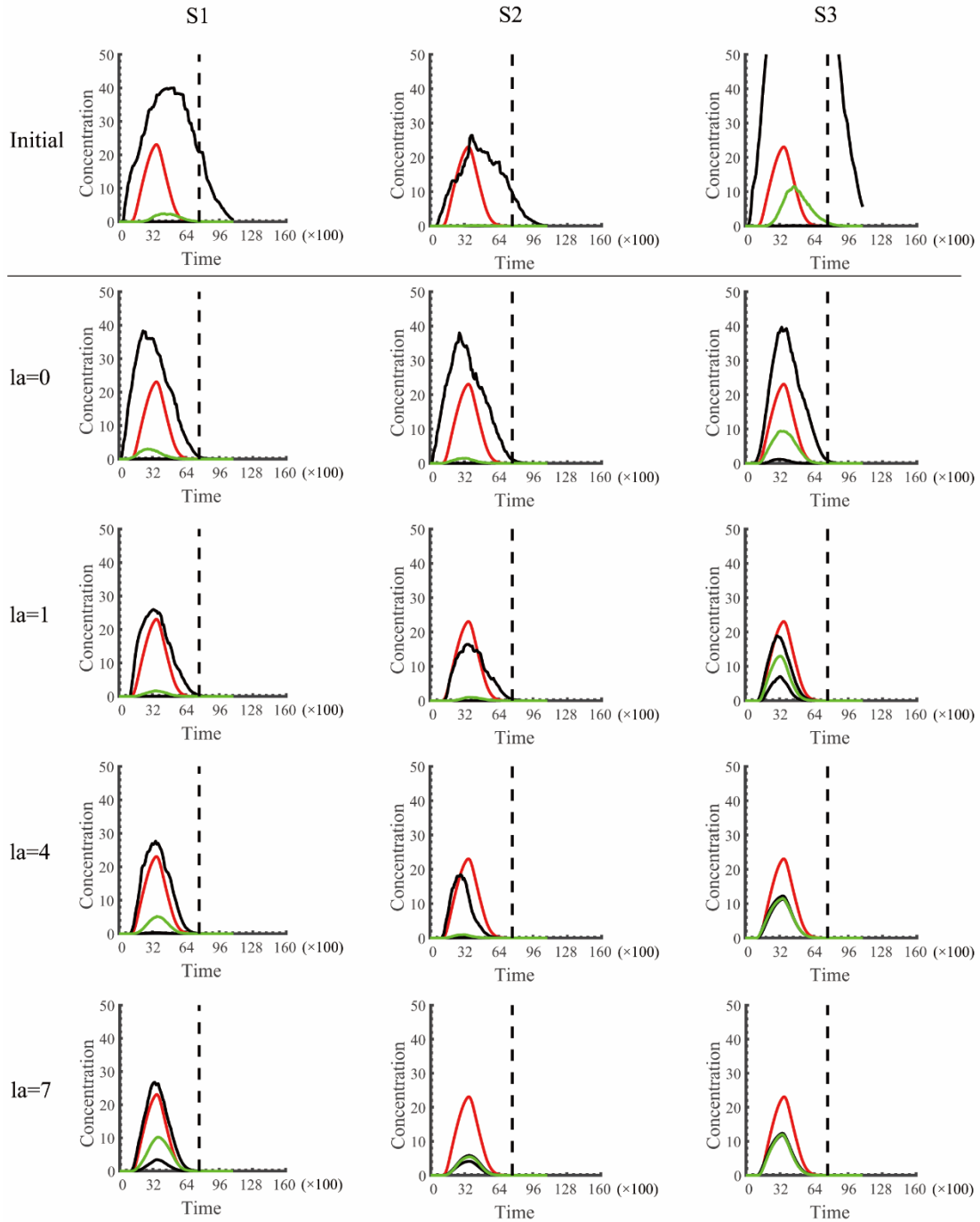
**Figure 10**: Boxplots computed with the initial and updated ensembles of release information parameters of contaminant source, including $T^0$, $\Delta T$ and $M$, after the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ data assimilation iterations in scenarios S1-S3. The dashed horizontal black line corresponds to the reference value.

To verify the accuracy of irregular source identification, two validation wells (#1, #2) are selected to predict the concentration evolution over time for the three scenarios, as shown in Figures 11 and 12. The initial concentrations exhibits significant uncertainties within the 5% to 95%, which decrease as assimilation iterations are performed. In S2 and S3, the confidence intervals nearly overlap but deviate from the true evolution. Particularly in well #2, the significant deviations indicate an underestimation of concentration uncertainty. In contrast, the reference evolutions for S1 consistently fall within the confidence intervals, demonstrating superior performance of ES-MDA-Bs with adequately equipped knot points in predicting concentration.
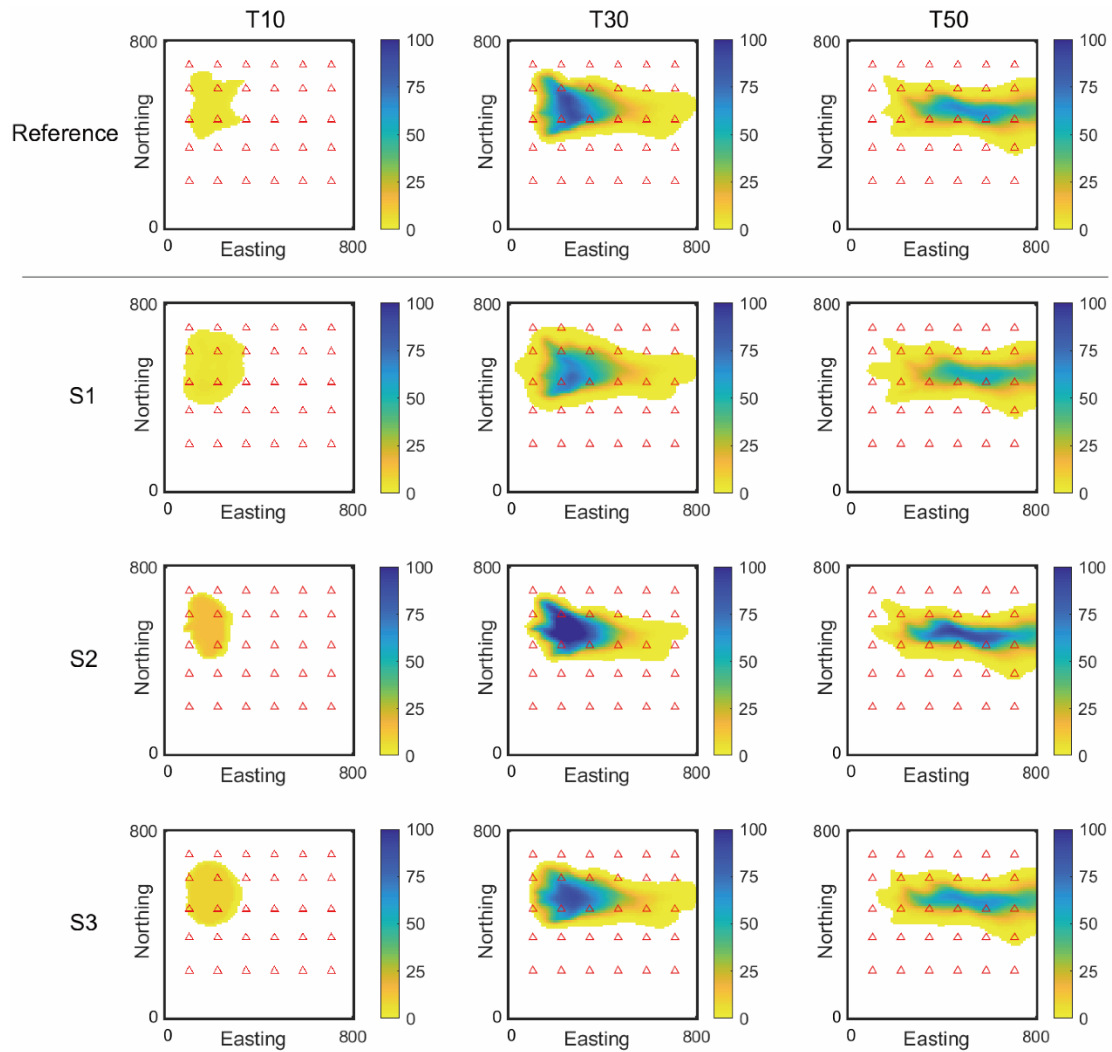
**Figure 11**: Time evolution of solute concentrations at the verification well #1 computed with the initial and updated ensembles of source parameters after the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ data assimilation in scenarios S1-S3. The left column corresponds to S1, the center column corresponds to S2 and the right column corresponds to S3. The red line corresponds to the reference field. The black lines correspond to the 5 and 95 percentiles of all realizations, and the green line corresponds to the median. The vertical dashed lines mark the end of the assimilation period.

472

**Figure 12**: Time evolution of solute concentrations at the verification well #2 computed with the initial and updated ensembles of source parameters after the $0^{th}$, $1^{st}$, $4^{th}$, and $7^{th}$ data assimilation in scenarios S1-S3. The left column corresponds to S1, the center column corresponds to S2 and the right column corresponds to S3. The red line corresponds to the reference field. The black lines correspond to the 5 and 95 percentiles of all realizations, and the green line corresponds to the median. The vertical dashed lines mark the end of the assimilation period.

Figure 13 presents a comparative analysis between the reference and the ensemble mean of the contaminant plume for the three scenarios at the $10^{th}$, $30^{th}$, and $50^{th}$ time steps. The
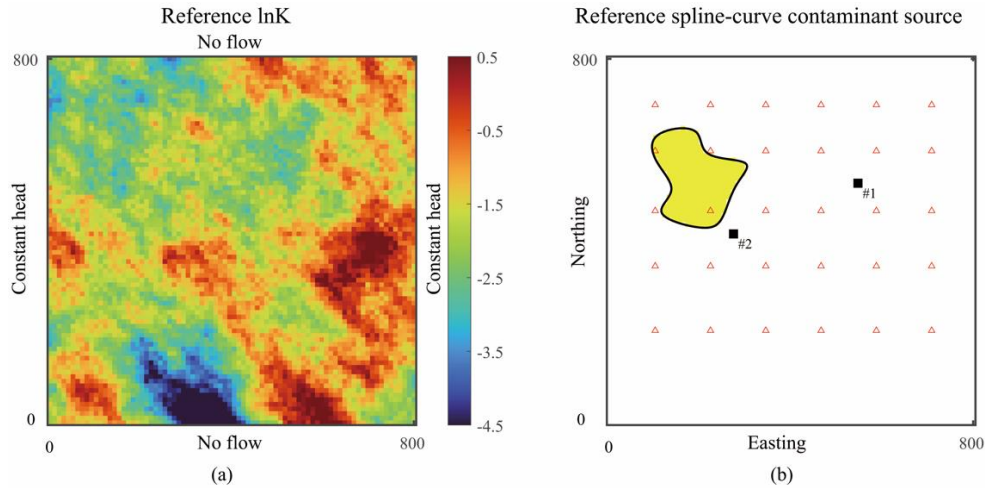
482 concentration plume is predicted using the updated source parameters from the 7th iteration.

483 The results show that S2 significantly over-predicts the concentration, while the S1 and S3 fit

484 well with the reference. However, S3 fails to capture the source irregularities derived from

485 the underlying source parameters, as previously discussed. This discrepancy may be attributed

486 to the fact that the predicted plume shape is influenced by multiple factors, such as

487 uncertainties in source location, spatial structure, and release history, as well as the accuracy

488 of these input parameters. These findings further confirm the accuracy of ES-MDA-Bs with

489 appropriately configured knot points to reproduce the contaminant plume.

490

491



492 **Figure 13**: Contaminant plume at the 10th, 30th, and 50th simulation time steps, computed with

493 the reference and the updated parameters after the 7th assimilation iteration in scenarios S1-

494 S3. From top to bottom: plume in reference field; ensemble mean of plumes calculated by

495 updates in S1; ensemble mean of plumes calculated by updates in S2; ensemble mean of

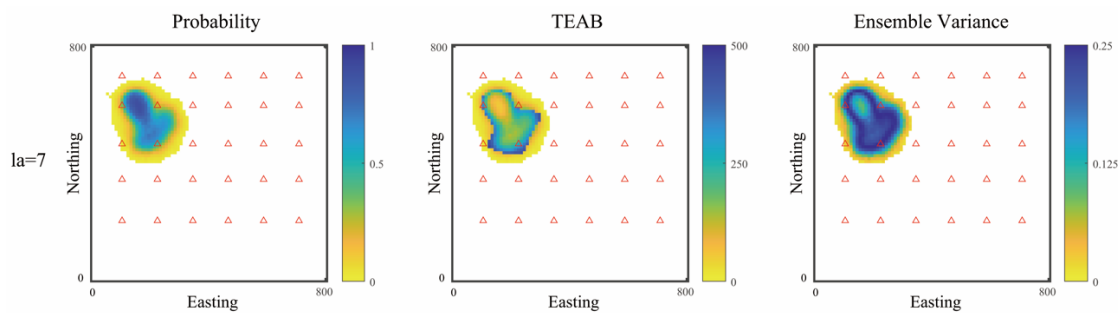496    plumes calculated by updates in S3.

497    To further evaluate the robustness and generalizability of the ES-MDA-Bs framework

498    configured with 10 knot points for identifying irregular non-point sources, a new synthetic

499    confined aquifer is designed, characterized by distinct hydrogeological properties and

500    contaminant source configurations compared to scenarios S1-S3 presented in Section 3.

501    Figure 14(a) illustrates the natural logarithm of a heterogeneous hydraulic conductivity field,

502    generated as a multi-Gaussian random field using the parameters specified in Table 1. In the

503    groundwater flow simulation, constant head boundaries are assigned to the eastern and

504    western edges of the aquifer, set at 90 [L] and 280 [L], respectively. For solute transport, the

505    longitudinal dispersivity is held constant in both space and time, while the ratio of horizontal

506    transverse to longitudinal dispersivity is reduced from 0.5 to 0.4, reflecting a potential

507    reduction in lateral dispersion due to altered velocity gradients induced by the modified

508    hydraulic conductivity distribution. As noted, a uniquely shaped non-point source, shown in

509    Figure 14(b), is introduced in this setup. The plausible range for each basic source parameter

510    is detailed in Table 3.

511



(a)                                                (b)

513    **Figure 14**: Reference log-conductivity field (a) and reference contamination source (b) in the
514    new synthetic aquifer. An indication of the type of boundary conditions used for the solution
515    of the flow equation is also shown in (a). The red triangles are the observation locations and
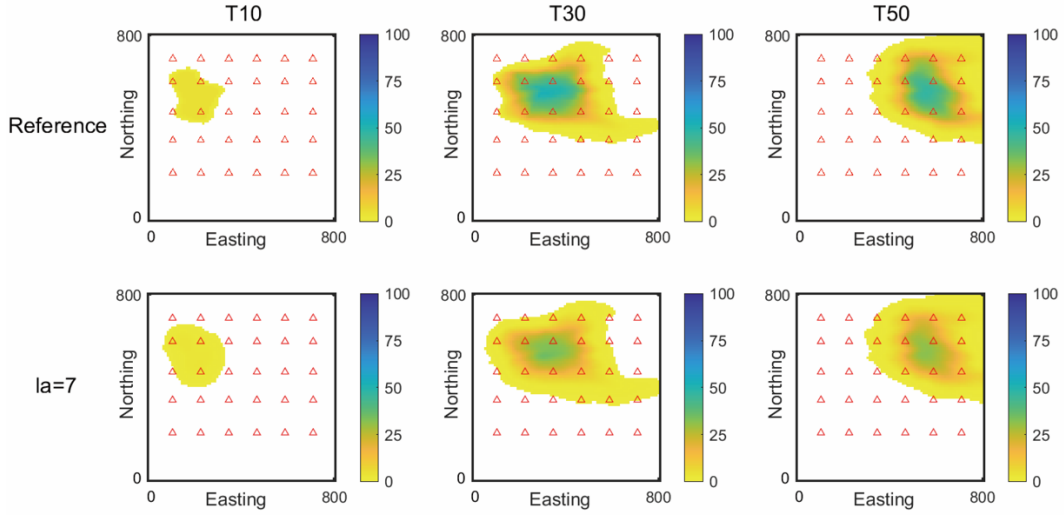516    the black squares are the validation locations shown in (b).

517    Figure 15 presents three statistical metrics—probability, $TEAB$, and ensemble variance

518     of the potential source—computed using the updated ensembles of positional parameters

519     following the 7th data assimilation iteration. A distinctly irregular spatial configuration is

520     evident in the updated spatial extent, characterized by pronounced ensemble variability along

521     the spline boundaries and localized positional discrepancies predominantly situated along

522     high-curvature regions in the southern boundary—both features are consistent with expected

523     behavior. These results further confirm that, as previously observed, ES-MDA-Bs with 10

524     knot points effectively effectively mitigates uncertainty underestimation while maintaining a

525     high degree of accuracy in spatial reconstruction, exhibiting only minor deviations in

526     proximity to the source margins.

527



528

529 **Figure 15**: Probability, TEAB, and ensemble variance of the potential source, calculated
530 using the updated ensembles of positional parameters following the 7th data assimilation
531 iteration in the new synthetic aquifer.

532     Figure 16 illustrates a comparative analysis between the reference contaminant plume

533     and the ensemble mean plume predicted using the updated source parameters from the 7th

534     iteration, conducted in the new synthetic aquifer at the 10th, 30th, and 50th time steps. The

535     ensemble prediction exhibits strong agreement with the reference plume in terms of temporal

536     evolution, spatial extent, and concentration levels. These findings further validate the

537     robustness of the ES-MDA-Bs method, configured with 10 knot points, in accurately

538     identifying randomly shaped, irregular non-point sources—including their location, structural

539     characteristics, and release history—and in reliably predicting contaminant plume distribution

540     across diverse synthetic scenarios.

541

**Figure 16**: Contaminant plume at the 10th, 30th, and 50th simulation time steps, computed with the reference and the updated parameters after the 7th assimilation iteration in the new synthetic aquifer. From top to bottom: plume in reference field; ensemble mean of plumes calculated by updates in the new synthetic aquifer.

However, transitioning the method from synthetic cases to real-world applications represents a significant challenge. In the next phase, we will adhere to the established research trajectory commonly adopted in groundwater point source identification studies (e.g., Xu and Gómez-Hernández, 2016; Chen et al., 2023b), progressing systematically from theoretical validation to sandbox experiments, and ultimately to large-scale field implementations.

## 5. Discussion

The aforementioned results demonstrate that the proposed ES-MDA-Bs effectively reproduces the complex spatial structure and the release information of a non-point source. However, this study represents an initial step in exploring irregular non-point sources. To support real-world applications, several key aspects still need further consideration:

(1) Performance evaluation and algorithm enhancement for heterogeneous release: In this study, ES-MDA-Bs is applied to identify an irregular non-point source with homogeneous releases. However, in real-world settings, source releases typically exhibit spatial heterogeneity and temporal fluctuations. Therefore, our future research will focus on evaluating and enhancing the ES-MDA-Bs algorithm, such as by integrating localization

562 techniques, to accurately identify the locations, spatial distributions, and varying release

563 characteristics of complex irregular non-point sources.

564 (2) Computational efficiency improvement: Data assimilation algorithms face

565 computational challenges due to the time-consuming nature of multiple runs on forward

566 model. Implementing a suitable surrogate model can alleviate this burden while maintaining

567 accuracy in reproducing state variable fields. In future work, we will develop a new algorithm

568 that integrates ES-MDA-Bs with an enhanced surrogate model to efficiently and accurately

569 reconstruct irregular non-point sources.

570 (3) Impact of observation noise level: This study adopts a fixed observation noise level

571 to isolate and evaluate the performance of ES-MDA-Bs in identifying irregular non-point

572 sources. However, given the intrinsic characteristics of ES-MDA algorithms, an

573 underestimated noise level may result in excessively large Kalman gains, leading to

574 overfitting of spurious residuals and potential ensemble collapse during sequential updates.

575 Therefore, future work will investigate the robustness of ES-MDA-Bs under varying noise

576 levels and develop adaptive strategies to mitigate errors induced by noise misspecification.

577 (4) Optimization of observation well placement: Practical constraints—such as

578 geological heterogeneity and economic limitations—often restrict the number and spatial

579 configuration of observation wells in field applications. Future research could focus on

580 overcoming these challenges by developing a multi-objective optimization framework for

581 well placement, integrated with inverse modeling techniques (e.g., ES-MDA-Bs), to

582 efficiently address complex non-point source identification problems under cost constraints.


583 **6. Summary**

584 In this paper, we evaluate the performance of ES-MDA-Bs in identifying a complex

585 irregular non-point source. Our findings demonstrate that ES-MDA-Bs effectively determines

586 the location, spatial structure and release information of the source, offering improved

587  accuracy over ES-MDA-RE. With the updated source parameters, the contaminant

588  concentration plume is successfully reproduced through forward prediction.

589      We further demonstrate that well-configured knot points significantly enhance the

590  accuracy of source reconstruction and effectively mitigates uncertainty underestimation.

591  Conversely, limited-configured knot points induce filter inbreeding, which can be mitigated

592  by employing the localization technique. Besides, we accurately estimate the initial release

593  time and release duration, though a slight overestimation of release mass-loading is observed.

594  Moreover, ES-MDA-Bs outperforms ES-MDA-RE in approximating the spatial distribution

595  of irregular sources, which is restricted to identify ellipse-like sources.

596      Building upon the research conducted by Xu et al. (2022), this paper extends the study

597  by proposing a new algorithm, ES-MDA-Bs, to accurately reproduce the location, spatial

598  structure, and release information of a complex irregular non-point source. This research has

599  made progress toward real-world applications. For future work, we intend to explore which

600  surrogate model combined with ES-MDA-Bs can improve the computational efficiency of

601  identifying complex irregular non-point sources. In addition, further testing and refining the

602  algorithm's performance to accommodate changes in source release patterns will be a priority.

609  **References**

610  An, Y., Y. Zhang, and X. Yan (2022), An integrated bayesian and machine learning approach

611      application to identification of groundwater contamination source parameters, Water, 14

612      (15), 2447. doi: 10.3390/w14152447.

Anshuman, A., and T. Eldho (2022), Entity aware sequence to sequence learning using lstms for estimation of groundwater contamination release history and transport parameters, Journal of Hydrology, 608, 127,662. doi: 10.1016/j.jhydrol.2022.127662.

Anshuman, A., and T. Eldho (2023), A parallel work ow framework using encoder-decoder lstms for uncertainty quantification in contaminant source identification in groundwater, Journal of Hydrology, 619, 129,296. doi: 10.1016/j.jhydrol.2023.129296.

Ayvaz, M. T. (2016), A hybrid simulation-optimization approach for solving the areal groundwater pollution source identification problems, Journal of Hydrology, 538, 161-176. doi: 10.1016/j.jhydrol.2016.04.008.

Bai, T., and P. Tahmasebi (2022), Characterization of groundwater contamination: A transformer-based deep learning model, Advances in Water Resources, 164, 104,217. doi: 10.1016/j.advwatres.2022.104217.

Butera, I., and M. G. Tanda (2003), A geostatistical approach to recover the release history of groundwater pollutants, Water Resources Research, 39 (12). doi: 10.1029/2003WR002314.

Chang, Z., W. Lu, and Z. Wang (2022), Study on source identification and source-sink relationship of lnapls pollution in groundwater by the adaptive cyclic improved iterative process and monte carlo stochastic simulation, Journal of Hydrology, 612, 128,109. doi: 10.1016/j.jhydrol.2022.128109.

Chang, Z., Guo, Z., Chen, K., et al. (2024). A comparison of inversion methods for surrogate‑based groundwater contamination source identification with varying degrees of model complexity. Water Resources Research, 60(4). doi: 10.1029/2023WR036051.

Chen, Z., J. J. Gómez-Hernández, T. Xu, and A. Zanini (2018), Joint identification of contaminant source and aquifer geometry in a sandbox experiment with the restart ensemble kalman filter, Journal of hydrology, 564, 1074-1084. doi: 10.1016/j.jhydrol.2018.07.073.

Chen, Z., L. Zong, J. J. Gómez-Hernández, T. Xu, Y. Jiang, Q. Zhou, H. Yang, Z. Jia, and S. Mei (2023a), Contaminant source and aquifer characterization: An application of es-mda demonstrating the assimilation of geophysical data, Advances in Water Resources, p. 104555. doi: 10.1016/j.advwatres.2023.104555.

643    Chen, Z., T. Xu, J. J. Gómez-Hernández, A. Zanini, and Q. Zhou (2023b), Reconstructing
644        the release history of a contaminant source with different precision via the ensemble
645        smoother with multiple data assimilation, Journal of Contaminant Hydrology, 252,
646        104,115. doi: 10.1016/j.jconhyd.2022.104115.

647    de Boor, C. (1978), A Practical Guide to Splines, Springer-Verlag.

648    De Boor, C., and C. De Boor (1978), A practical guide to splines, vol. 27, springerverlag
649        New York.

650    Emerick, A. A., and A. C. Reynolds (2013), Ensemble smoother with multiple data
651        assimilation, Computers & Geosciences, 55, 3-15. doi: 10.1016/j.cageo.2012.03.011.

652    Evensen, G., and P. J. Van Leeuwen (2000), An ensemble kalman smoother for non-linear
653        dynamics, Monthly Weather Review, 128 (6), 1852-1867. doi: 10.1175/1520-
654        0493(2000)1282.0.CO;2.

655    Farin, G. (2014), Curves and surfaces for computer-aided geometric design: a practical
656        guide, Elsevier.

657    Ge, Y., W. Lu, and Z. Pan (2023), Groundwater contamination source identification based
658        on sobol sequences-based sparrow search algorithm with a bilstm surrogate model,
659        Environmental Science and Pollution Research, 30 (18), 53,191-53,203. doi:
660        10.1007/s11356-023-25890-0.

661    Gómez-Hernández, J. J., and A. G. Journel (1993), Joint sequential simulation of
662        multigaussian fields, in Geostatistics Tróia'92: Volume 1, pp. 85-94, Springer. doi:
663        10.1007/978-94-011-1739-5_8.

664    Gorelick, S. M., B. Evans, and I. Remson (1983), Identifying sources of groundwater
665        pollution: An optimization approach, Water Resources Research, 19 (3), 779-790. doi:
666        10.1029/WR019i003p00779.

667    Gzyl, G., A. Zanini, R. Fraczek, and K. Kura (2014), Contaminant source and release
668        history identification in groundwater: a multi-step approach, Journal of contaminant
669        hydrology, 157, 59-72. doi: 10.1016/j.jconhyd.2013.11.006.

670    Hou, Z., W. Lao, Y. Wang, and W. Lu (2021), Homotopy-based hyper-heuristic searching
671        approach for reciprocal feedback inversion of groundwater contamination source and

672 aquifer parameters, Applied Soft Computing, 104, 107,191. doi:

673 10.1016/j.asoc.2021.107191.

674 Hwang, H.-T., S.-W. Jeen, D. Kaown, S.-S. Lee, E. A. Sudicky, D. T. Steinmoeller, and K.-

675 K. Lee (2020), Backward probability model for identifying multiple contaminant source

676 zones under transient variably saturated flow conditions, Water Resources Research, 56

677 (4), e2019WR025,400. doi: 10.1029/2019WR025400

678 Hwang, J. C., and R. M. Koerner (1983), Groundwater pollution source identification from

679 limited monitoring well data: Part 1|theory and feasibility, Journal of hazardous materials,

680 8 (2), 105-119. doi: 10.1016/0304-3894(83)80050-8.

681 Jamshidi, A., J. M. V. Samani, H. M. V. Samani, A. Zanini, M. G. Tanda, and M. Mazaheri

682 (2020), Solving inverse problems of unknown contaminant source in groundwater-river

683 integrated systems using a surrogate transport model based optimization, Water, 12 (9),

684 2415. doi: 10.3390/w12092415.

685 Jha, M., and B. Datta (2015), Application of dedicated monitoring{network design for

686 unknown pollutant-source identification based on dynamic time warping, Journal of

687 Water Resources Planning and Management, 141 (11), 04015,022. doi:

688 10.1061/(ASCE)WR.1943-5452.0000513.

689 Jin, X., G. Mahinthakumar, E. M. Zechman, and R. S. Ranjithan (2009), A genetic

690 algorithm-based procedure for 3d source identification at the borden emplacement site,

691 Journal of Hydroinformatics, 11 (1), 51-64. doi: 10.2166/hydro.2009.002.

692 Juhász, I., and M. Hoffmann (2004), Constrained shape modification of cubic b-spline

693 curves by means of knots, Computer-Aided Design, 36 (5), 437-445. doi: 10.1016/S0010-

694 4485(03)00116-7.

695 Kontos, Y. N., T. Kassandros, K. Perifanos, M. Karampasis, K. L. Katsifarakis, and K.

696 Karatzas (2022), Machine learning for groundwater pollution source identification and

697 monitoring network optimization, Neural Computing and Applications, 34 (22), 19,515-

698 19,545. doi: 10.1007/s00521-022-07507-8.

699 Li, J., Z. Wu, H. He, and W. Lu (2023a), Identifying groundwater contamination sources

700 based on the hybrid grey wolf gradient algorithm and deep belief neural network,

701 Stochastic Environmental Research and Risk Assessment, 37 (5), 1697-1715. doi:

702     10.1007/s00477-022-02360-6.

703     Li, L., H. Zhou, H.-J. Hendricks Franssen, and J. J. Gómez-Hernández (2012), Modeling

704         transient groundwater flow by coupling ensemble kalman filtering and upscaling, Water

705         Resources Research, 48 (1). doi: 10.1029/2010WR010214.

706     Li, Y., W. Lu, Z. Pan, Z. Wang, and G. Dong (2023b), Simultaneous identification of

707         groundwater contaminant source and hydraulic parameters based on multilayer

708         perceptron and flying foxes optimization, Environmental Science and Pollution

709         Research, pp. 1-15. doi: 10.1007/s11356-023-27574-1.

710     Luo, C., W. Lu, Z. Pan, Y. Bai, and G. Dong (2023), Simultaneous identification of

711         groundwater pollution source and important hydrogeological parameters considering the

712         noise uncertainty of observational data, Environmental Science and Pollution Research,

713         pp. 1-16. doi: 10.1007/s11356-023-28091-x.

714     Lyche, T., and K. Mørken (1999), The sensitivity of a spline function to perturbations of the

715         knots, BIT Numerical Mathematics, 39, 305-322. doi: 10.1023/A:1022346030560.

716     Mahar, P. S., and B. Datta (1997), Optimal monitoring network and ground-water pollution

717         source identification, Journal of water resources planning and management, 123 (4), 199-

718         207. doi: 10.1061/(ASCE)0733-9496(1997)123:4(199).

719     Mahinthakumar, G., and M. Sayeed (2005), Hybrid genetic algorithm-local search methods

720         for solving groundwater source identification inverse problems, Journal of water

721         resources planning and management, 131 (1), 45-57. doi: 10.1061/(ASCE)0733-

722         9496(2005)131:1(45).

723     Mirghani, B. Y., K. G. Mahinthakumar, M. E. Tryby, R. S. Ranjithan, and E. M. Zechman

724         (2009), A parallel evolutionary strategy based simulation-optimization approach for

725         solving groundwater source identification problems, Advances in Water Resources, 32

726         (9), 1373-1385. doi: 10.1016/j.advwatres.2009.06.001.

727     Mirghani, B. Y., E. M. Zechman, R. S. Ranjithan, and G. Mahinthakumar (2012),

728         Enhanced simulation-optimization approach using surrogate modeling for solving

729         inverse problems, Environmental Forensics, 13 (4), 348-363. doi:

730         10.1080/15275922.2012.702333.

731     Pan, Z., W. Lu, Z. Chang, et al. (2021), Simultaneous identification of groundwater

732     pollution source spatial-temporal characteristics and hydraulic parameters based on deep

733     regularization neural network-hybrid heuristic algorithm, Journal of Hydrology, 600, 126,

734     586. doi: 10.1016/J.JHYDROL.2021.126586.

735 Pan, Z., W. Lu, and Y. Bai (2022), Groundwater contamination source estimation based on

736     a refined particle filter associated with a deep residual neural network surrogate,

737     Hydrogeology Journal, 30 (3), 881-897. doi: 10.1007/s10040-022-02454-z.

738 Pan, Z., W. Lu, H. Wang, and Y. Bai (2023a), Groundwater contaminant source

739     identification based on an ensemble learning search framework associated with an auto

740     xgboost surrogate, Environmental Modelling & Software, 159, 105,588. doi:

741     10.1016/j.envsoft.2022.105588.

742 Pan, Z., W. Lu, and Y. Bai (2023b), Groundwater contaminated source estimation based on

743     adaptive correction iterative ensemble smoother with an auto lightgbm surrogate, Journal

744     of Hydrology, 620, 129,502. doi: 10.1016/j.jhydrol.2023.129502.

745 Pan, Z., Guo, Z., Chen, K., et al. (2025). A deep adaptive bidirectional generative

746     adversarial neural network (Bi-GAN) for groundwater contamination source estimation.

747     Journal of Hydrology, 132753. doi: 10.1016/j.jhydrol.2025.132753.

748 Park, H., and J.-H. Lee (2007), B-spline curve fitting based on adaptive curve refinement

749     using dominant points, Computer-Aided Design, 39 (6), 439-451. doi:

750     10.1016/j.cad.2006.12.006.

751 Piegl, L., and W. Tiller (1996), The NURBS book, Springer Science & Business Media.

752 Schoenberg, I. J. (1946), Contributions to the problem of approximation of equidistant data

753     by analytic functions. part b. on the problem of osculatory interpolation. a second class of

754     analytic approximation formulae, Quarterly of Applied Mathematics, 4 (2), 112-141.

755 Singh, P., Mahor, V., Lakshmaiya, N., Shanker, K., Kaliappan, S., & Muthukannan, M., et

756     al. (2024). Prediction of groundwater contamination in an open landfill area using a novel

757     hybrid clustering-based ai model. Environment Protection Engineering, 50(1). doi:

758     10.37190/epe240106.

759 Srivastava, D., and R. M. Singh (2015), Groundwater system modeling for simultaneous

760     identification of pollution sources and parameters with uncertainty characterization,

761     Water resources management, 29, 4607-4627. doi: 10.1007/s11269-015-1078-8.

762 Wu, Y., Li, M., Xie, H., et al. (2025). Characterizing multi-source heavy metal

763 contaminated sites at the Hun River basin: An integrated deep learning and data

764 assimilation approach. Journal of Hydrology, 648, 132349. doi:

765 10.1016/j.jhydrol.2024.132349.

766 Xu, T., and J. J. Gómez-Hernández (2016), Joint identification of contaminant source

767 location, initial release time, and initial solute concentration in an aquifer via ensemble

768 kalman filtering, Water Resources Research, 52 (8), 6587-6595. doi:

769 10.1002/2016WR019111.

770 Xu, T., and J. J. Gómez-Hernández (2018), Simultaneous identification of a contaminant

771 source and hydraulic conductivity via the restart normal-score ensemble kalman filter,

772 Advances in Water Resources, 112, 106-123. doi: 10.1016/j.advwatres.2017.12.011.

773 Xu, T., J. J. Gómez-Hernández, H. Zhou, and L. Li (2013), The power of transient

774 piezometric head data in inverse modeling: An application of the localized normal-score

775 enkf with covariance in ation in a heterogenous bimodal hydraulic conductivity field,

776 Advances in Water Resources, 54, 100-118. doi: 10.1016/j.advwatres.2013.01.006.

777 Xu, T., J. J. Gómez-Hernández, Z. Chen, and C. Lu (2021), A comparison between es-mda

778 and restart enkf for the purpose of the simultaneous identification of a contaminant source

779 and hydraulic conductivity, Journal of Hydrology, 595,125,681. doi:

780 10.1016/j.jhydrol.2020.125681.

781 Xu, T., W. Zhang, J. J. Gómez-Hernández, Y. Xie, J. Yang, Z. Chen, and C. Lu (2022),

782 Non-point contaminant source identification in an aquifer using the ensemble smoother

783 with multiple data assimilation, Journal of Hydrology, 606, 127,405. doi:

784 10.1016/j.jhydrol.2021.127405.

785 Zhang, X., Jiang, S., Zheng, N., et al. (2024a). Integration of DDPM and ILUES for

786 simultaneous identification of contaminant source parameters and non-Gaussian

787 channelized hydraulic conductivity field. Water Resources Research, 60(9),

788 e2023WR036893. doi: 10.1029/2023WR036893.

789 Zhang. W., T. Xu, Z. Chen, J. J. Gómez-Hernández, C. Lu, J. Yang, Y. Ye, and M. Jing

790 (2024b), Simultaneous identification of a non-point contaminant source with gaussian

791 spatially distributed release and heterogeneous hydraulic conductivity in an aquifer using

792 the les-mda method, Journal of Hydrology, 630, 130,745.

793    doi:10.1016/j.jhydrol.2020.125681.

794    Zheng, N., Li, Z., Xia, X., et al. (2024). Estimating line contaminant sources in non-

795    Gaussian groundwater conductivity fields using deep learning-based framework. Journal

796    of Hydrology, 630, 130727. doi: 10.1016/j.jhydrol.2024.130727.