

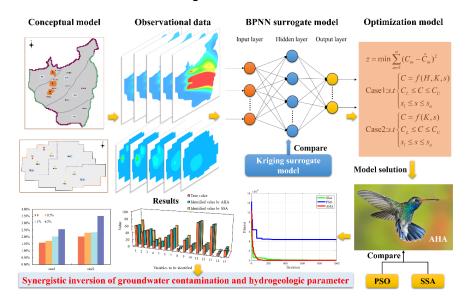


Synergistic identification of hydrogeological parameters and pollution 1 source information for groundwater point and areal source 2 contamination based on machine learning surrogate-artificial 3 hummingbird algorithm 4 Chengming Luo¹, Xihua Wang^{1,2*}, Y. Jun Xu³, Shunqing Jia¹, Zejun Liu¹, Boyang Mao¹, 5 Qinya Lv¹, Xuming Ji¹, Yanxin Rong¹, Yan Dai¹ 6 7 ¹ College of Civil Engineering, Tongji University, 1239 Siping Road, Shanghai 200092, 8 China9 ² Department of Earth and Environmental Sciences, University of Waterloo, ON N2L 10 3G1, Canada 11 ³ School of Renewable Natural Resources, Louisiana State University Agricultural 12 Center, Baton Rouge, Louisiana, USA 13 14 Email: 21531@tongji.edu.cn (Xihua Wang) 15 Tel.: + 0086 0431 13089410676; Fax: + 0086 021 65986809 16 17 18 19





21 Graphical Abstract



22 23





24	Highlights
25	A highly adaptable inversion framework is adapted to different groundwater pollution
26	scenarios.
27	Synergetic identification of source information, hydraulic conductivity and boundary
28	condition in PSC.
29	The artificial hummingbird algorithm is applied to solve the optimized model.





Abstract

30

31 Effectively remediating groundwater contamination relies on the precise determination of its sources. In recent years, a growing research focus has been placed on concurrently 32 33 estimating hydrogeological characteristics and locating pollutant origins. However, the 34 identification of precise synergistic identification of point and areal contamination sources of groundwater and combined hydrogeological parameters has not been 35 36 effectively solved. This study developed an inversion framework that integrates 37 machine learning surrogates with the artificial hummingbird algorithm (AHA). The 38 surrogate models approximating the simulation system were constructed using both backpropagation neural networks (BPNN) and Kriging techniques. The AHA was then 39 employed to solve the optimized model, and its performance was benchmarked against 40 particle swarm optimization (PSO) and the sparrow search algorithm (SSA). The 41 42 applicability of this inversion framework was assessed by application to point sources of contamination (PSC) and areal source contamination (ASC). The robustness of the 43 framework was verified through application to scenarios with different noise levels. 44 45 The results showed that surrogate model constructed by the BPNN method provided estimates that were closer to those of the simulation model in comparison to the kriging 46 method, coefficient of determination (R²) is 0.9994 and mean relative error (MRE) is 47 3.70% in PSC, and R² is 0.9989 and MRE is 4.48% in ASC. The performance of the 48 AHA exceeded those of the PSO and the SSA. In PSC, MRE of the identification result 49 is 1.58%; In ASC, MRE of the identification result is 2.03%, with the AHA able to 50 rapidly and accurately identify the global optimum and improve the inversion efficiency. 51

https://doi.org/10.5194/egusphere-2025-2083 Preprint. Discussion started: 20 May 2025 © Author(s) 2025. CC BY 4.0 License.





- 52 The proposed inversion framework was demonstrated to apply to both groundwater
- 53 PSC and ASC problems with strong robustness, providing a reliable basis for
- 54 groundwater pollution remediation and management.
- 55 **Keywords:** Groundwater contamination identification; Synergistic identification; Point
- and areal sources contamination; Surrogate model; Artificial hummingbird algorithm





1 Introduction

57

58 Groundwater pollution adversely affects human production and life (Wang et al., 2022; Liu et al., 2024). The remediation of groundwater contamination is important for 59 ensuring human health and socioeconomic development. However, groundwater 60 61 contamination is difficult to detect and treat due to its hidden nature, thereby complicating the assessment of groundwater pollution risk and contamination liability 62 63 (Li et al., 2021). Remediation requires the identification of sources of groundwater 64 contamination (location, number, release history, etc.) and hydrogeological conditions 65 (Daranond et al., 2020; Pan et al., 2022b). However, directly obtaining this information can pose a challenge, with a proven method being the identification of groundwater 66 contamination by inversion of limited observational data. 67 68 Inversion of groundwater aquifer hydrogeologic parameters and pollution source 69 information is a widely studied topic. In past studies on groundwater contamination identification (GCI), many researchers have focused on the separate identification of 70 hydrogeological parameters or pollution source information. For example, Singh and 71 72 Datta (2007) utilized backpropagation-based artificial neural network techniques specifically for the identification of groundwater pollution sources. Similarly, Mahar 73 and Datta (2000) employed a nonlinear optimization model to identify the location, 74 duration, and magnitude of the contamination source. Liu et al. (2022) inverted 75 76 hydrogeological parameters through a simulation-optimization approach, while Wang et al. (2024a) combined three different inversion algorithms and a kriging surrogate 77 model to invert hydraulic conductivity. While simplifying the problem, these methods 78





allow researchers to focus on specific aspects. However, although the individual 79 80 identification method can be effective in some cases, it often overlooks the interconnectivity between hydrogeological parameters and pollution sources. 81 Currently, the simultaneous identification of hydrogeological parameters and 82 83 pollution source information is gaining increasing attention in research. Researchers have employed various advanced technologies to achieve this goal. Wang et al. (2021) 84 85 utilized a parallelized heuristic algorithm to concurrently determine both aquifer 86 characteristics and the groundwater pollution sources. Pan et al. (2021) integrated a 87 Bayesian-regularized deep neural network surrogate to jointly infer pollution source details and hydraulic conductivity. Hou et al. (2021) integrated homotopy-based inverse 88 optimization theory with a multi-kernel extreme learning machine to finish the co-89 90 identification of contamination sources and aquifer parameters. Luo et al. (2023) 91 leveraged machine learning techniques to establish an inverse relationship between model outputs and inputs, enabling fast and simultaneous retrieval of pollution source 92 attributes and hydrogeological properties. Although these methods have advanced the 93 94 field, improving recognition accuracy remains a major challenge in the simultaneous 95 identification process. The simulation-optimization method has been widely applied in GCI research 96 because of its robust mathematical foundation (Mirghani et al., 2009) and its ability to 97 98 identify multiple variables simultaneously. To enhance both identification accuracy and efficiency using simulation-optimization, two key approaches are employed: one is to 99 optimize the model solution method for better performance, and the other is to construct 100

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122





a surrogate model with high approximation accuracy. Optimizing the model solution method is essential. Since heuristic optimization algorithms are more capable of identifying global optima, many have been applied to GCI. Mirghani et al. (2012) implemented a genetic algorithm within optimization to identify sources of contamination. Jiang et al. (2013) combined a harmony search algorithm with a contamination transport simulation model to characterize contamination sources. Additional methods, such as simulated annealing (Rao, 2006; Yeh et al., 2007; Jha and Datta, 2013) and sparrow search algorithms (SSA) (Pan et al., 2022b), have also been applied to GCI. However, increasing dimensionality and complexity in GCI problems make it difficult for many optimization algorithms to efficiently search for global optima. Constructing high-accuracy surrogate models is another crucial strategy. Surrogate models can significantly reduce computation time and improve inversion efficiency. Among these models, the widely used kriging (Chugh et al., 2018; Zhang et al., 2019; Jiang et al., 2020) and backpropagation neural network (BPNN) (Sargolzaei et al., 2012; Zhang et al., 2021; Wang et al., 2024b) methods offer high flexibility and strong nonlinear fitting capabilities. Despite these advances, previous studies have overly focused on point source contamination (PSC) or areal source contamination (ASC) scenarios in isolation. However, the identification of precise synergistic identification of PSC and ASC of groundwater and combined hydrogeological parameters has not been effectively solved. Based on the above problems, this paper proposes an inversion framework integrating a machine learning surrogate model with the artificial hummingbird





algorithm (AHA) using the simulation-optimization method (Fig. 1). Both BPNN and 123 124 kriging were utilized to develop surrogate models for the simulation model. AHA was introduced to solve the optimization model, with its solution results compared against 125 those of PSO and SSA. The applicability of this inversion framework was evaluated 126 127 through its application to both PSC and ASC scenarios. The objectives of this study were: (1) To assess the performance of the surrogate models constructed by BPNN and 128 129 kriging, and to identify the model with better practicality and adaptability to replace the 130 simulation model; (2) To examine the advantages of AHA for solving the optimization 131 model by comparing it with other optimization algorithms under the same conditions, and to further enhance the model's solving accuracy; (3) To apply the simulation-132 optimization-based inversion framework to complete the inversion tasks for PSC and 133 ASC scenarios and validate the framework's effectiveness, while also evaluating its 134 135 robustness through tests under different noise conditions.

2. Methodology

136

137

142

143

144

2.1. Simulation model

In this study, the numerical groundwater simulation framework comprised both a flow component and a solute transport module. The fundamental two-dimensional (2D) partial differential equation governing groundwater flow is formulated as follows:

141
$$\frac{\partial}{\partial x_i} (K_{ij} (H - z) \frac{\partial H}{\partial x_j}) + W = \mu \frac{\partial H}{\partial t} (x, y) \in S \ i, j \in 1, 2 \ t \ge 0$$
 (1)

where K_{ij} is hydraulic conductivity, W is the volumetric flux per unit volume, μ is the specific yield, H is the water level elevation, z is the elevation of the aquifer floor, and S is the boundary of the spatial domain.

153

154

155

156

157

158





$$\frac{\partial C}{\partial t} = \frac{\partial}{\partial x_i} (D_{ij} \frac{\partial C}{\partial x_j}) - \frac{\partial}{\partial x_i} (u_i C) + \frac{R}{n_e}$$
 (2)

$$u_{i} = \frac{K_{ij}}{n_{e}} \frac{\partial H}{\partial x_{i}}$$
 (3)

where C denotes the contaminant concentration in groundwater, t is the temporal

variable, u_i indicates the average flow velocity, R accounts for source and sink contributions, D_{ij} refers to the hydrodynamic dispersion tensor, and n_e represents the effective porosity of the medium. To obtain numerical solutions for the groundwater

151 flow and solute transport equations, the MODFLOW and MT3DMS packages were

employed (Asher et al., 2015).

2.2. Kriging method

Kriging was employed to develop the underlying framework of the approach by capturing both the correlation and stochastic variability of variables within a confined spatial domain, thereby enabling the estimation of optimal regional values. The association between input and output variables is described through a regression-based expression as shown below (Zhao et al., 2022a):

159
$$y(x) = \sum_{i=1}^{k} \beta_{1i} f_i(x) + z(x)$$
 (4)

where $\hat{y}(x)$ is the estimated value of pollutant concentration y(x), $f_i(x)(i =$

161 $1, \dots, k$) is the basis function of the known regression model, and z(x) is the random

162 part.

163 The following equations were satisfied:

164
$$\begin{cases} E(z(x)) = 0 \\ D(z(x)) = \sigma^2 \\ cov[z(x_i), z(x_j)] = \sigma^2 R(x_i, x_j) \end{cases}$$
 (5)





- where $R(x_i, x_i)$ is the correlation function between the sampled point x_i and x_i .
- 166 $(i = 1, 2, \dots, m; j = 1, 2, \dots, m)$
- The Gaussian model is commonly used:

168
$$R(x_{i}, x_{j}) = \exp\left(-\sum_{k=1}^{m} \theta_{k} \left| x_{k_{i}} - x_{k_{j}} \right|^{2}\right)$$
 (6)

- where θ_k is a coefficient to be determined, which can be obtained by calculation.
- 170 2.3. The BPNN method
- 171 A typical back-propagation neural network (BPNN) is composed of three
- fundamental components (Fig. 2): (1) an input layer, (2) the hidden layers, and (3) an
- 173 output layer. The computation process proceeds in two main phases: forward
- propagation and backward propagation (Chen et al., 2010; Zhang et al., 2018).
- 175 1) During forward propagation, data are introduced into the network via the input
- 176 layer, and subsequently processed through successive layers to yield the final output.
- 177 BPNNs frequently employ a nonlinear sigmoid activation function:

178
$$f(x) = \frac{1}{1 + e^{-x}} \tag{7}$$

The calculation of the forward transmission output layer is:

180
$$I_{j} = \sum_{i=1} w_{ij} o_{i} + b \qquad o_{j} = f(I_{j}) = \frac{1}{1 + e^{I_{j}}}$$
 (8)

- where O_i represents the output of neuron i, O_j is the output of neuron j, b is the bias
- term, and W_{ij} is the weight of the connection between neuron i and neuron j.
- 2) Backward propagation involves the random assignment of the weight of the first
- positive feedback process within the output layer. The adjustment of the parameters of
- the entire network is required. Network adjustment is performed by minimizing the
- discrepancy between the predicted output and the target category in the output layer.





187 Specifically, for the output layer:

188
$$E_{i} = O_{i}(1 - O_{i})(T_{i} - O_{i})$$
 (9)

- where E_j represents the error value at the jth node and T_j denotes the corresponding
- output. The hidden layer's output is determined by summing the weighted contributions
- 191 from the errors of the lower nodes:

192
$$E_{i} = O_{i}(1 - O_{i}) \sum_{k} E_{k} W_{ik}$$
 (10)

- where E_k is the error gradient for the subsequent node k and W_{jk} is the weight connecting
- 194 node j to t node k. Following error calculation, the weight is adjusted according to the
- 195 error gradient:

$$\Delta W_{ij} = \eta E_j O_i$$

$$W_{ij} = W_{ij} + \Delta W_{ij}$$
(11)

197 where η is the learning rate.

198 2.4 Artificial Hummingbird Algorithm (AHA)

199 The AHA consists of three main elements: food sources, hummingbirds, and the visit 200 table. Hummingbirds typically assess food sources based on factors such as nectar quality, individual flower nectar content, and replenishment rates. For simplicity, it can 201 be assumed that all food sources share the same flower type and number. 202 203 Hummingbirds within a population can exchange information, be assigned to specific 204 food sources, track nectar replenishment rates, and record the duration each food source remains unvisited. The visit table records the time since a hummingbird last visited a 205 food source, and is used to assign visit levels; hummingbirds can harvest more nectar 206 207 by first accessing food sources with higher access levels, following which food sources with the highest nectar replenishment rate are chosen (Zhao et al., 2022b). The AHA is 208





- 209 algorithmically described below.
- 210 (1) Initialization
- Firstly, n humming birds are randomly placed on n food sources:

$$x_i = Low + r \cdot (Up - Low) \quad i = 1, \dots, n$$
 (12)

213 The access table for the food source is then initialized:

214
$$VT_{i,j} = \begin{cases} 0 & \text{if } i \neq j \\ \text{null } i = j \end{cases} i = 1, ..., n; j = 1, ..., n$$
 (13)

- 215 where Low and Up are the lower and upper boundaries for a d-dimensional problem
- 216 respectively, r represents a random vector of [0,1], and x_i is the position of the ith food
- source. For i = j, $VT_{i,j} = null$ indicates the sourcing of food from a specific source.
- For $i \neq j$, $VT_{i,j} = 0$ indicates that the ith humming bird has just visited the jth food
- 219 source in the current iteration.
- 220 (2) Guided foraging
- 221 Hummingbirds identify food sources in two steps: (1) identifying the food source
- 222 with the highest access level; (2) selecting the food source with the highest nectar
- 223 replenishment rate. After identifying the target food source, the hummingbird can fly to
- 224 the target source to feed. During foraging, direction switching vectors used to control
- 225 the availability of one or more directions in the D-dimensional space are introduced to
- 226 model three flight skills: omnidirectional, diagonal, and axial flight. These flight
- 227 models can be extended to the d-D space, and the mathematical model of axial flight is:





228
$$D^{(i)} = \begin{cases} 1 & if \quad i = randi([1, d]) \\ 0 & else \end{cases} i = 1, ..., d$$
 (14)

Diagonal flight is defined as:

230
$$D^{(i)} = \begin{cases} 1 & if \quad i = P(j), j \in [1, k] \\ P = randperm(k), k \in [2, [r_1 \cdot (d-2)] + 1] & i = 1, ..., d \\ 0 & else \end{cases}$$
 (15)

231 Omnidirectional flight is defined as:

232
$$D^{(i)} = 1 \quad i = 1, ..., d$$
 (16)

- where randi([1,d]) is a randomly generated integer from 1 to d, randperm(k)
- creates a random permutation of integers from 1 to k, and r_1 is a random number in the
- 235 range of 0 to 1.
- 236 Hummingbirds can access and obtain target food sources through these flight abilities.
- 237 New food sources identified during the search are recorded along with previously
- 238 identified food sources. The guided foraging behavior and candidate food sources can
- 239 be represented as:

240
$$v_i(t+1) = x_{i tar}(t) + a \cdot D \cdot (x_i(t) - x_{i tar}(t))$$
 (17)

$$241 a \sim N(0,1) (18)$$

- 242 where $x_{i,tar}(t)$ is the location of the food source that the *i*th hummingbird plans to
- visit, $x_i(t)$ represents the location of the *i*th food source at time t, and a is a leading
- 244 factor obeying a normal distribution.
- The location of the *i*th food source is updated as:





- 247 where $f(\cdot)$ represents the function fitness value. The formula for updating the location
- 248 can contribute to the preferential selection of food sources with a high nectar supply
- 249 rate.
- 250 (3) Territorial foraging
- Since the quality of food sources within a foraging area may vary, hummingbirds
- 252 actively search within that area. The regional foraging strategies and candidate food
- 253 sources of hummingbirds can be represented as:

254
$$v_i(t+1) = x_i(t) + b \cdot D \cdot x_i(t)$$
 (20)

$$255 b \sim N(0,1) (21)$$

- 256 where b is a territorial factor obeying a normal distribution. Eq. (20) allows different
- 257 hummingbirds to use their specific flight skills to identify new food sources near the
- 258 target source.
- 259 (4) Migration foraging
- 260 Migration coefficients are defined in the AHA algorithm to prevent the generation of
- 261 local optimums. The exceedance of the number of iterations of the set migration
- 262 coefficient results in the hummingbird located in the worst food source repeating a
- search for a new food source across the entire search range and the subsequent updating
- of the visit table.





$$x_{wor}(t+1) = Low + r \cdot (Up - Low) \tag{22}$$

where x_{wor} is the food source with the worst nectar supply rate. The migration 266 coefficient relative to population size can be defined as.

$$M = 2n \tag{23}$$

3. Case studies 269

267

270

271

272

273

274

275

276

277

281

282

283

284

285

The present study designed a groundwater PSC case study and an ASC case study to verify the applicability of the proposed GCI framework. Since the present study established two hypothetical examples, a set of variables to be identified and background variables for input into the groundwater contamination simulation model were established for each example for forward computation. The pollutant concentrations monitored at wells were used as observed data. The robustness of the inversion framework was verified by adding random noise to the observed data, expressed as:

278
$$\alpha_1 = \alpha(1 + l \cdot \text{rand}), \ l = 0.5\%, 1\% \text{ and } 2\%$$
 (24)

279 where α represents the observation data, α_1 indicates observation data with added 280 noise, l is the max disturbance range, and rand is a random number between -1 and 1.

3.1 Case study 1: groundwater PSC

The study area is 2,500 m and 1,400 m from east to west and north to south, respectively, with topography decreasing from west to east and groundwater flow from northwest to southeast. The study area contains an inhomogeneous isotropic aquifer, and the present study focused on a layer of diving aquifer with a thickness of 10 m (Table 1).

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307





Groundwater flow was represented as 2D steady flow, and the study area was divided into three areas according to differences in hydraulic conductivities. Since the northern and southern parts of the study area are very weakly permeable formations, they were generalized in the present study as no-flow boundaries. Rivers formed the boundaries of the western and eastern parts, and were generalized as specific head boundaries (Fig. 3). In this case study, the variables to be identified fell into three main categories: (1) head values at the specific head boundaries, including H_1 and H_2 ; (2) hydraulic conductivities for each part of the study area, including K_1 , K_2 , and K_3 ; (3) the intensities of the release of pollutants from the two sources during the release periods: $S = S_a T_b$; a = 1, 2; and b = 1, 2, 3, 4, 5 (Table 3). S_aT_b represents the intensity of pollution source aduring the bth stress period; this case study had a study period of 10 years (Table 1, Fig. 4), with both sources only releasing pollutants in the first five years (Table 2). Five wells were established to monitor the concentrations of groundwater contaminants once a year. The study area was spatially discretized into 50 m \times 50 m grids (Table 1). 3.2 Case study 2: groundwater ASC The present study selected the hypothetical case study used by Pan et al. (2022a) as a case study. The site has an area of 5 km², with a length of 2.5 km and width of 2 km from east to west and south to north, respectively. Groundwater flows from northwest to southeast. The study area was conceptualized as an inhomogeneous isotropic aquifer and the current study focused on a diving aquifer, in which flow was represented as 2D steady flow. The study area's aquifers were categorized into four zones based on





hydraulic conductivity, labeled K_1 to K_4 . The western and eastern river boundaries were 308 modeled as specified head boundaries, while the northern and southern regions, 309 characterized by low permeability granite, were treated as no-flow boundaries (Fig. 5, 310 Table 4). 311 312 Within this case study, the variables to be identified fell into two categories: (1) hydraulic conductivities of each part of the study area, including K_1 to K_4 ; (2) the 313 intensities of pollutants released by three areal sources of contamination: $S = S_a T_b$; a =314 315 1, 2, 3; and b = 1, 2, 3, 4, 5 (Table 5). $S_a T_b$ indicates the intensity of pollution source a 316 during the bth stress period. A total of nine monitoring wells were established to monitor the concentrations of groundwater contaminants once a year (Fig. 6). The study area 317 was spatially discretized as 20 m × 20 m grids (Table 4). 318 319 4. Model construction 320 4.1 Establishment of surrogate models The present study established two case studies: the PSC and the ASC. The variables to 321 be identified for the PSC case study included three categories with 15 dimensions, 322 323 whereas those to be identified for the ASC case study included two categories with 19 dimensions. The present study used the Latin hypercube method to sample within the 324 feasible domain of the variables to be identified. This sampling process was 325 implemented in MATLAB. Sample groups for the PSC and ASC case studies totaled 326 390 and 490, respectively, and the input sample dataset was generated by random 327 combination. 328 The parameters obtained from the above sampling were input into the groundwater 329





simulation model. The simulation model was then run to obtain the pollutant concentrations at the 390 and 490 monitoring groups in the PSC and ASC case studies, respectively. These simulated pollutant concentrations were used as the output sample dataset, and the output sample dataset was combined with the input sample dataset to form the input-output sample dataset. The kriging and BPNN methods were used to establish the surrogate models of the simulation model. The first 350 and 440 groups of the PSC and ASC case input-output sample datasets, respectively, were used as training samples in each case study to construct surrogate models, while the remaining 40 and 50 groups were used as test samples to evaluate the accuracy of the surrogate models.

The present study applied the coefficient of determination (R²), the mean relative error (MRE), and the root mean square error (RMSE) to assess the accuracy of the fit of the estimations of the surrogate models to the output of the simulation model.

1) R²: The closer R² to 1, the more accurate the surrogate model is.

344
$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \overline{y}_{i})^{2}}$$
 (25)

2) MRE: The average deviation between the outputs of the surrogate model and the outputs of the simulation model.

$$MRE = \frac{\sum_{i=1}^{n} \left| \frac{y_i - \hat{y}_i}{y_i} \right|}{n}$$
 (26)

348 3) RMSE: The value of the RMSE is inversely proportional to the fitting accuracy of the surrogate model.





350
$$RMSE = \sqrt{\frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n}}$$
 (27)

where \overline{y}_i is the average true value, n is the number of samples, \hat{y}_i is the output of the surrogate model, y_i is the true value of the variable to be identified.

4.2 Establishment of the optimization models

This study employed the CGI through the S-O method, which consists of two main components: a groundwater contaminant transport simulation model and an optimization model aimed at minimizing the least squares error between the simulated and true values. To reduce the computational burden caused by repeated simulation calls, a surrogate model was used in place of the simulation model. While the same objective function was applied in both case studies, there were minor variations in the decision variables and constraints. The decision variables chosen for case study 1 included the boundary head values, the hydraulic conductivities of the site, and the release history of the contaminant source; those for case study 2 included the hydraulic conductivities of the site and the release history of the contaminant source. The constraint conditions were influenced by the decision variables. The optimization was expressed as:

$$z = \min \sum_{m=1}^{n} (C_m - \hat{C}_m)^2$$

$$Case1:s.t \begin{cases} C = f(H, K, s) \\ C_L \le C \le C_U \\ s_l \le s \le s_u \end{cases}$$

$$Case2:s.t \begin{cases} C = f(K, s) \\ C_L \le C \le C_U \\ s_l \le s \le s_u \end{cases}$$

$$(28)$$





where z is the objective function, C_m is the monitored pollutant concentration in the 367 368 mth monitoring well, \hat{C}_m is the simulated pollutant concentration in the mth monitoring well, C is the pollutant concentration, H is the head value at the boundary, 369 370 s is the pollution source intensity, k represents the hydraulic conductivities of the site, 371 C_L and C_U are the upper and lower bound values of pollutant concentration, respectively, and s_l and s_u are the upper and lower bound values of pollution source 372 373 intensity, respectively. 374 The AHA was used to identify the optimal combination of parameters according to 375 the objective function through multiple iterative calculations, with this parameter set adopted as the result of inversion. The numbers of hummingbird populations and 376 iterations were set to 500 and 1,000, respectively. 377 378 5. Results 379 5.1 Surrogate models The surrogate model for case study 1 using the kriging method achieved an R² of 0.9942, 380 MRE of 13.43%, and RMSE of 11.8262 (Table 6), while the BPNN method produced 381 382 values of 0.9994, 3.70%, and 3.6526, respectively (Table 6). Similarly, for case study 2, the kriging method yielded an R² of 0.9837, MRE of 9.98%, and RMSE of 37.7547, 383 whereas the BPNN method provided corresponding values of 0.9989, 3.70%, and 384 3.6526 (Table 6). The BPNN method demonstrated superior goodness-of-fit statistics 385 386 compared to the kriging method in both case studies. While the simulation model required 500 hours for 1,000 iterations, the BPNN 387 surrogate model completed the same number of iterations in 67 seconds, significantly 388





reducing the computation time.

5.2 Optimization algorithms

The BPNN surrogate model was embedded into the optimization model to optimize the parameter combination according to the objective function. This study employed AHA within the optimization process and compared its performance against SSA and PSO under the same population size and number of iterations. In the optimization of case study 1, PSO failed to converge after reaching the maximum number of iterations, while AHA and SSA converged after 120 and 350 iterations, respectively (Fig. 7a). For case study 2, both PSO and SSA failed to converge within the maximum number of iterations, whereas AHA converged after 150 iterations (Fig. 7b).

Given the results from case study 1, where both AHA and SSA converged, the subsequent analysis focused on these two algorithms. AHA achieved an optimal search value closer to the true value and reached the global optimum, while SSA settled at a local optimum (Fig. 8). These results demonstrate that AHA not only converged faster than SSA but also identified the global optimum, thereby improving the accuracy and efficiency of GCI.

5.3 Inversion results and robustness assessment

The BPNN-AHA inversion framework developed in this study was applied to identify groundwater PSC and ASC and obtain inversion values. To verify the framework's robustness and reliability, random noise levels of 0.5%, 1%, and 2% were added to the observed data. The average relative errors under each noise level were recorded (Table 7, Table 8). The highest inversion accuracy was achieved in the noise-free case for both





case study 1 and case study 2, with average relative errors of 1.58% and 2.03%, respectively (Table 9). At a 0.5% noise level, the average relative errors for case study 1 and case study 2 were 1.71% and 2.3%. At 1% noise, they were 2.03% and 2.33%, while at 2% noise, they increased to 2.55% and 3.52%, respectively. Although noise impacted the inversion accuracy, the framework maintained high performance, with the average relative errors for both case studies remaining below 5% (Fig. 9). These results confirm the strong robustness and stability of the proposed inversion framework.

6 Discussion

6.1 Analysis of surrogate models

In the current research on GCI, it is common to use machine learning methods to construct surrogate models for groundwater simulation. Various methods are employed, such as long short-term memory neural networks (Li et al., 2021), light gradient boosting machines (Pan et al., 2023), and deep residual networks (Xu et al., 2024b), each with its own advantages. This study focuses on adaptable methods. Compared to the methods mentioned above, the BPNN surrogate model developed in this paper features a simple structure, high flexibility, and broad adaptability. It performs well in different scenarios, including the PSC and ASC cases analyzed in this paper, where the R² values are 0.9994 and 0.9989 and the MRE values are 3.7% and 4.48%, respectively. These results demonstrate the model's excellent ability to fit the input-output relationship of the simulation model. The effectiveness of a surrogate model lies not only in its complexity but also in how well it fits the problem at hand. A good surrogate model should maintain both high accuracy and strong adaptability. In this paper, the





ASC is drawn from Pan et al. (2022a), which had been widely validated in other studies. 433 434 For example, Li et al. (2023) used the same case to validate an inversion method, applying a multilayer perceptron model to the simulation, achieving an R² of 0.9999 435 and an MRE of 2.85%. Similarly, Xu et al. (2024a) employed automatic machine 436 437 learning methods for surrogate model construction, achieving an R2 of 0.9754 and an MRE of 4.154%. Compared to the surrogate models developed by these researchers, 438 439 the BPNN model constructed in this study also demonstrates excellent approximation 440 accuracy, further validating the advantages of the proposed method. 441 6.2 Analysis of optimization algorithms This paper compares the AHA with PSO and SSA under the same preconditions and 442 finds that AHA offers clear advantages in both convergence speed and global 443 444 optimization capability. Based on these results, AHA was chosen to solve the optimization model, and its adaptability was further verified in two different cases. In 445 the field of optimization algorithms, the "no free lunch principle" (Zhao et al., 2022b) 446 emphasizes that no single algorithm performs well across all optimization problems. 447 448 When addressing real-world problems, it is essential to understand the nature of the problem thoroughly before selecting the appropriate optimization algorithm. This 449 principle encourages researchers to develop new and more effective algorithms from 450 different perspectives, providing more options for optimization problem researchers. 451 452 This insight also applies to groundwater pollution traceability. Given the diverse nature of pollution traceability problems, it is challenging for any single optimization 453 algorithm to be universally applicable. As research deepens, these problems tend to 454





become more high-dimensional and nonlinear, necessitating the exploration of algorithms with stronger global optimization capabilities and higher search efficiency. Additionally, it is important to consider alternative uses of optimization methods. One promising approach involves using optimization techniques to improve machine learning models by identifying optimal parameters (hyperparameters) during training, which can significantly enhance model accuracy (Jia et al., 2024).

6.3 Inversion analysis

Previous studies related to GCI employed a variety of methods to conduct either single or simultaneous inversion characterization of pollution sources and to identify hydrogeological parameters of the model. Li et al. (2022) identified the number, location, and release history of pollution sources, while Li et al. (2008) focused on determining the hydraulic conductivities of a study site. Bai et al. (2022) utilized inversion techniques to simultaneously characterize pollution sources and identify the hydraulic conductivities within their simulation models. While some studies have applied inversion to the boundary conditions of the simulation model (Jiao et al., 2019), fewer studies have simultaneously characterized pollution sources and identified both hydrogeological parameters and boundary conditions of the model. Source information, model hydrogeological parameters, and boundary conditions are all critical components of groundwater contamination simulation models. Inaccuracies in any of these components can affect the overall results of inversion, making it essential to identify all components simultaneously. Therefore, in the PSC case of this study, the release history of the pollutant source, the hydraulic conductivity of the model, and the specific head





boundary values were simultaneously identified. This simultaneous identification of 477 478 multiple key parameters enhances the reliability and effectiveness of decision support 479 systems. The overall inversion framework in this paper combines BPNN and AHA and is 480 481 validated under different noise scenarios to account for the effect of noise in the observed data. The results indicate that the inversion framework demonstrates high 482 483 robustness. However, a limitation of this paper is that noise is not addressed, and its 484 presence can contaminate the observed data, further impacting the accuracy of GCI. 485 Noise elimination methods could be applied to the observed data in future studies. Another major limitation is the generalization of the actual groundwater system. 486 Groundwater systems are often complex, necessitating model simplifications through 487 assumptions (e.g., homogeneity, isotropy) that may not reflect the actual geological 488 conditions, thereby affecting model accuracy. To address actual problems, the 489 hydrogeological conditions of the study area should be thoroughly investigated, 490 ensuring the model closely represents the actual situation, reducing error, improving 491 492 model accuracy, and ultimately enhancing inversion accuracy. 7 Conclusions 493 In this study, a BPNN-AHA inversion framework was developed to accurately and 494 synergistically identify groundwater point and areal sources of contamination and 495 496 combined hydrogeologic parameters. Among them, the BPNN surrogate model can well replace the simulation model, and the AHA had good global optimization 497 capability and excellent solution accuracy. The robustness of the proposed methodology 498





499 was verified by applying the inversion framework to scenarios with different noise 500 levels. The conclusions of the present study are listed below: (1) The construction of a surrogate model to the simulation model satisfied the fitting 501 accuracy requirement while also significantly reducing the computational time. The 502 503 current study established BPNN and kriging surrogate models, with a comparison of the outputs of the models illustrating that the former obtained a higher fitting accuracy, 504 505 leading to its application in the inversion framework. 506 (2) The present study applied AHA within the model optimization, with the results 507 compared to those of PSO and SSA optimization. Compared to PSO and SSA, AHA rapidly reached convergence and identified the global optimum, thereby significantly 508 improving the accuracy and efficiency of inversion. 509 510 (3) The proposed inversion framework can realize the synergistic identification of PSC 511 and ASC combined with hydrogeological parameters, which can ensure high identification accuracy, and the inversion framework has strong robustness under 512 different noise levels. While individual identification simplifies the problem but may 513 514 ignore correlations between parameters, synergistic identification improves the accuracy and consistency of identification by synchronizing the estimation of pollution 515 sources and hydrogeological parameters. However, noise and parameter estimation 516 uncertainties may still affect the reliability of the inversion results. Therefore, 517 518 uncertainty analysis needs to be further considered in subsequent studies. Overall, the BPNN-AHA inversion framework has excellent inversion performance and strong 519 practicability, which can provide a reliable basis for groundwater pollution remediation 520

https://doi.org/10.5194/egusphere-2025-2083 Preprint. Discussion started: 20 May 2025 © Author(s) 2025. CC BY 4.0 License.





521 and management.





522 Ethical Approval The study did not use any data which need approval. Consent to Participate All authors participated in the process of draft completion. All 523 authors have read and agreed to the published version of the manuscript. 524 Consent to Publish All authors agree to publish. 525 526 **Author's contributions** Chengming Luo: Methodology, Formal analysis, Software, Conceptualization, 527 528 Validation, Writing-original draft. Xihua Wang: Supervision, Resources, Funding 529 acquisition, Writing- review & editing. Y. Jun Xu: Supervision, Writing- review & 530 editing. Shunqing Jia, Zejun Liu, & Boyang Mao: Software, Formal analysis. Qinya 531 Lv, Xuming Ji, Yanxin Rong & Yan Dai: Methodology, Conceptualization. **Declaration of competing interest** 532 533 The authors declared that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper. 534 Acknowledgment 535 This research was supported by the Fundamental Research Funds for the Central 536 537 Universities (02002150257) and Overseas High-level Talents Program of Shanghai and Leading Talents (Overseas) Program of Shanghai. 538 **Data Availability** 539 The code used in this study can be found at https://doi.org/10.5281/zenodo.14568110. 540 541

References

542





543 Asher, M.J., Croke, B.F.W., Jakeman, A.J., Peeters, L.J.M., 2015. A review of surrogate models and their application to groundwater modeling. Water Resources Research, 544 51(8): 5957-5973. 545 546 Bai, Y., Lu, W., Li, J., Chang, Z. and Wang, H., 2022. Groundwater contamination source identification using improved differential evolution Markov chain 547 548 algorithm. Environmental Science and Pollution Research, 29(13): 19679-19692. 549 Chen, D., Lu, J., Shen, Y., 2010. Artificial neural network modelling of concentrations 550 of nitrogen, phosphorus and dissolved oxygen in a non-point source polluted river in Zhejiang Province, southeast China. Hydrological Processes, 24(3): 290-299. 551 Chugh, T., Jin, Y., Miettinen, K., Hakanen, J., Sindhya, K., 2018. A Surrogate-Assisted 552 553 Reference Vector Guided Evolutionary Algorithm for Computationally Expensive 554 Many-Objective Optimization. Ieee Transactions on Evolutionary Computation, 22(1): 129-142. 555 Daranond, K., Yeh, T.-C.J., Hao, Y., Wen, J.-C., Wang, W., 2020. Identification of 556 557 groundwater basin shape and boundary using hydraulic tomography. Journal of Hydrology, 588. 558 Hou, Z., Lao, W., Wang, Y., Lu, W., 2021. Hybrid homotopy-PSO global searching 559 approach with multi-kernel extreme learning machine for efficient source 560 identification of DNAPL-polluted aquifer. Computers & Geosciences, 155. 561 Jha, M., Datta, B., 2013. Three-Dimensional Groundwater Contamination Source 562 Identification Using Adaptive Simulated Annealing. Journal of Hydrologic 563





564 Engineering, 18(3): 307-317. 565 Jiang, C. et al., 2020. Real-time estimation error-guided active learning Kriging method for time-dependent reliability analysis. Applied Mathematical Modelling, 77: 82-566 98. 567 568 Jiang, S., Zhang, Y., Wang, P., Zheng, M., 2013. An almost-parameter-free harmony search algorithm for groundwater pollution source identification. Water Science 569 570 and Technology, 68(11): 2359-2366. 571 Jia, S., Wang, X., Xu, Y.J., Liu, Z. and Mao, B., 2024. A New Data-Driven Model to 572 Predict Monthly Runoff at Watershed Scale: Insights from Deep Learning Method Applied in Data-Driven Model. Water Resources Management, 38(13): 5179-5194. 573 Jiao, J., Zhang, Y. and Wang, L., 2019. A new inverse method for contaminant source 574 575 identification under unknown solute transport boundary conditions. Journal of 576 Hydrology, 577. Li, J., Lu, W., Luo, J., 2021. Groundwater contamination sources identification based 577 on the Long-Short Term Memory network. Journal of Hydrology, 601. 578 579 Li, J., Lu, W. and Luo, J., 2021. Groundwater contamination sources identification based on the Long-Short Term Memory network. Journal of Hydrology, 601. 580 Li, J., Wu, Z., Lu, W. and He, H., 2022. Simultaneous identification of the number, 581 location and release intensity of groundwater contamination sources based on 582 583 simulation optimization and ensemble surrogate model. Water Supply, 22(10): 7671-7689. 584 Li, W., Englert, A., Cirpka, O.A. and Vereecken, H., 2008. Three-dimensional 585





geostatistical inversion of flowmeter and pumping test data. Ground Water, 46(2): 586 193-201. 587 Li, Y., Lu, W., Pan, Z., Wang, Z. and Dong, G., 2023. Simultaneous identification of 588 589 groundwater contaminant source and hydraulic parameters based on multilayer 590 perceptron and flying foxes optimization. Environmental Science and Pollution Research, 30(32): 78933-78947. 591 592 Liu, Y., Luo, J., Xiong, Y., Ji, Y. and Xin, X., 2022. Inversion of hydrogeological 593 parameters based on an adaptive dynamic surrogate model. Hydrogeology Journal, 594 30(5): 1513-1527. Liu, Z., Wang, X., Wan, X., Jia, S. and Mao, B., 2024. Evolution origin analysis and 595 health risk assessment of groundwater environment in a typical mining area: 596 597 Insights from water-rock interaction and anthropogenic activities. Environmental 598 Research, 252. Luo, C., Lu, W., Pan, Z., Bai, Y. and Dong, G., 2023. Simultaneous identification of 599 groundwater pollution source and important hydrogeological parameters 600 601 considering the noise uncertainty of observational data. Environmental Science and Pollution Research, 30(35): 84267-84282. 602 Mirghani, B.Y., Zechman, E.M., Ranjithan, R.S., Mahinthakumar, G., 2012. Enhanced 603 Simulation-Optimization Approach Using Surrogate Modeling for Solving Inverse 604 605 Problems. Environmental Forensics, 13(4): 348-363. Mahar, P.S. and Datta, B., 2000. Identification of pollution sources in transient 606 groundwater systems. Water Resources Management, 14(3): 209-227. 607





608	Mirghani, B.Y., Mahinthakumar, K.G., Tryby, M.E., Ranjithan, R.S. and Zechman,
609	E.M., 2009. A parallel evolutionary strategy based simulation-optimization
610	approach for solving groundwater source identification problems. Advances in
611	Water Resources, 32(9): 1373-1385.
612	Pan, Z., Lu, W., Bai, Y., 2022a. Groundwater contamination source estimation based on
613	a refined particle filter associated with a deep residual neural network surrogate.
614	Hydrogeology Journal, 30(3): 881-897.
615	Pan, Z., Lu, W., Wang, H., Bai, Y., 2022b. Recognition of a linear source contamination
616	based on a mixed-integer stacked chaos gate recurrent unit neural network-hybrid
617	sparrow search algorithm. Environmental Science and Pollution Research, 29(22):
618	33528-33543.
619	Pan, Z., Lu, W. and Bai, Y., 2023. Groundwater contaminated source estimation based
620	on adaptive correction iterative ensemble smoother with an auto lightgbm
621	surrogate. Journal of Hydrology, 620.
622	Pan, Z., Lu, W., Fan, Y. and Li, J., 2021. Identification of groundwater contamination
623	sources and hydraulic parameters based on bayesian regularization deep neural
624	network. Environmental Science and Pollution Research, 28(13): 16867-16879.
625	Rao, S.V.N., 2006. A computationally efficient technique for source identification
626	problems in three-dimensional aquifer systems using neural networks and
627	simulated annealing. Environmental Forensics, 7(3): 233-240.
628	Sargolzaei, J., Asl, M.H., Moghaddam, A.H., 2012. Membrane permeate flux and
629	rejection factor prediction using intelligent systems. Desalination, 284: 92-99.





630 Singh, R.M. and Datta, B., 2007. Artificial neural network modeling for identification 631 of unknown pollution sources in groundwater with partially missing concentration observation data. Water Resources Management, 21(3): 557-572. 632 Wang, H., Lu, W. and Chang, Z., 2021. Simultaneous identification of groundwater 633 634 contamination source and aquifer parameters with a new weighted-average wavelet variable-threshold denoising method. Environmental Science and 635 636 Pollution Research, 28(28): 38292-38307. 637 Wang, X., Xu, Y.J. and Zhang, L., 2022. Watershed scale spatiotemporal nitrogen 638 transport and source tracing using dual isotopes among surface water, sediments and groundwater in the Yiluo River Watershed, Middle of China. Science of the 639 Total Environment, 833. 640 Wang, Z., Yue, C. and Wang, J., 2024a. Evaluating parameter inversion efficiency in 641 642 Heterogeneous Groundwater models using Karhunen-Loève expansion: a comparative study of genetic algorithm, ensemble smoother, and MCMC. Earth 643 Science Informatics, 17(4): 3475-3491. 644 Wang, X., Ji, X., Xu, Y. J., Mao, B., Jia, S., & Wang, C., et al., 2024b. Multi-645 machine learning methods to predict spatial variation characteristics of total 646 nitrogen at watershed scale: evidences from the largest watershed (yangtze river 647 watershed), asian. Science of the Total Environment, 949. 648 649 Xu, Y. et al., 2024a. Groundwater contaminant source identification considering unknown boundary condition based on an automated machine learning surrogate. 650 Geoscience Frontiers, 15(1). 651





652	Xu, Y. et al., 2024b. Intelligent enhanced particle filter with deep residual network
653	surrogate for accurate groundwater pollution source characterization. Journal of
654	Hydrology, 642.
655	Yeh, HD., Chang, TH., Lin, YC., 2007. Groundwater contaminant source
656	identification by a hybrid heuristic approach. Water Resources Research, 43(9).
657	Zhang, P., Cai, Y., Wang, J., 2018. A simulation-based real-time control system for
658	reducing urban runoff pollution through a stormwater storage tank. Journal of
659	Cleaner Production, 183: 641-652.
660	Zhang, X., Wang, L., Sorensen, J.D., 2019. REIF: A novel active-learning function
661	toward adaptive Kriging surrogate models for structural reliability analysis.
662	Reliability Engineering & System Safety, 185: 440-454.
663	Zhang, Yg. et al., 2021. Application of an enhanced BP neural network model with
664	water cycle algorithm on landslide prediction. Stochastic Environmental Research
665	and Risk Assessment, 35(6): 1273-1291.
666	Zhao, Y., Fan, D., Li, Y., Yang, F., 2022a. Application of machine learning in predicting
667	the adsorption capacity of organic compounds onto biochar and resin.
668	Environmental Research, 208.
669	Zhao, W., Wang, L., Mirjalili, S., 2022b. Artificial hummingbird algorithm: A new bio-
670	inspired optimizer with its engineering applications. Computer Methods in
671	Applied Mechanics and Engineering, 388.

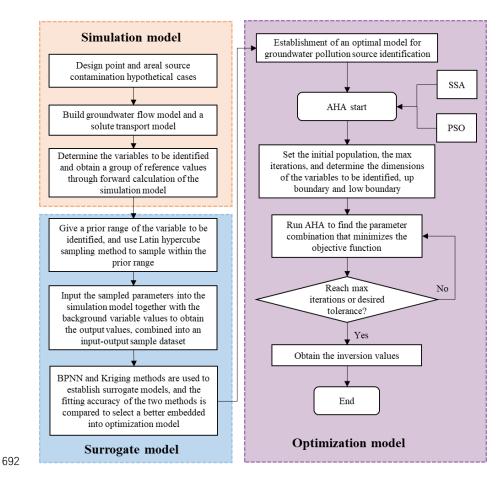




672 Figure captions Figure 1: General process used in the present study to construct the machine learning 673 surrogate model-artificial hummingbird algorithm framework. 674 Figure 2: Structure of a back-propagation neural network (BPNN). 675 676 **Figure 3:** Schematic diagram of case study 1. Figure 4: Distributions of concentrations of groundwater pollutants over different 677 678 periods: (a)–(j) represent 1–10 years. 679 **Figure 5:** Schematic diagram of case study 2. 680 Figure 6: Distributions of concentrations of groundwater pollutants over different periods: (a) 1 year; (b) 2 years; (c) 3 years; (d) 4 years; (e) 5 years. 681 Figure 7: Convergence curves of the sparrow search algorithm (SSA), particle swarm 682 683 optimization (PSO), and artificial hummingbird algorithm (AHA) applied to case study. (a) case study1; (b) case study2. 684 Figure 8: Comparison between the true values and optimal values for the sparrow 685 search algorithm (SSA) and artificial hummingbird algorithm (AHA). 686 Figure 9: Comparison of relative errors for case studies 1 and 2 under different noise 687 levels. 688 689 690 691

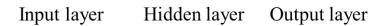


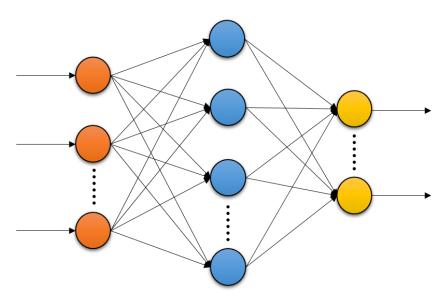






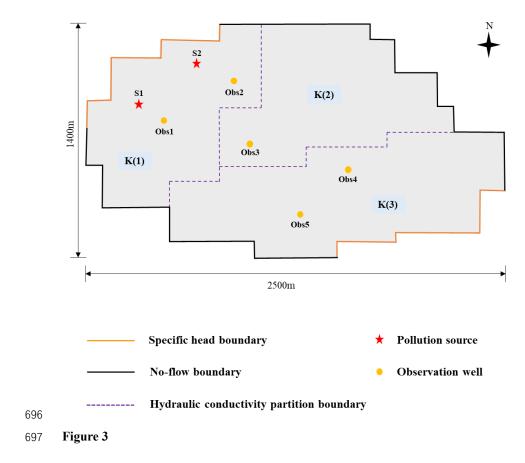






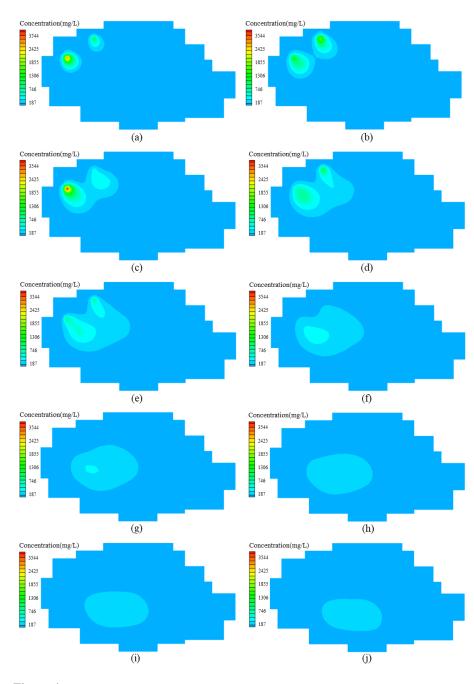








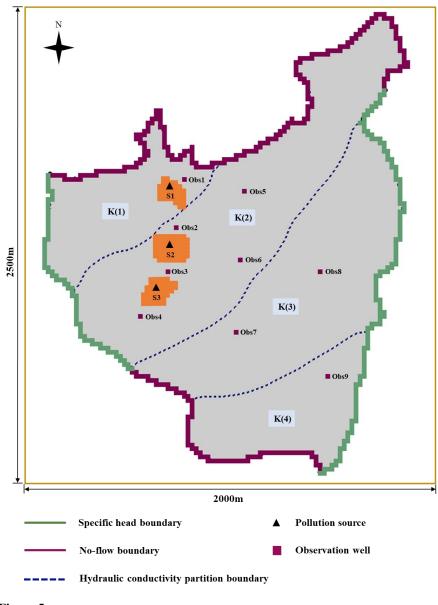




699 Figure 4

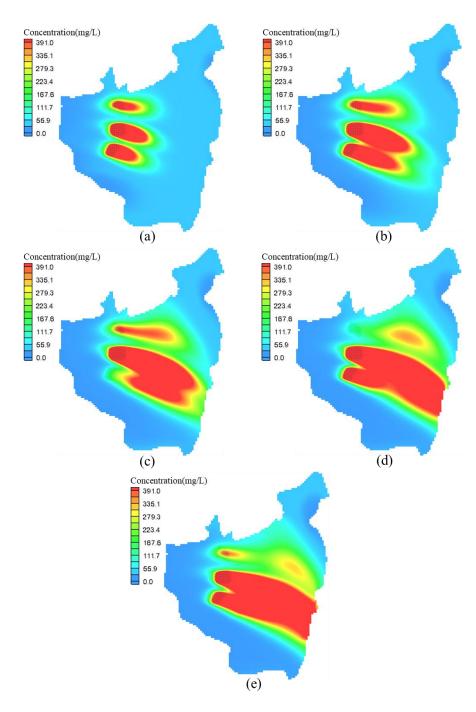








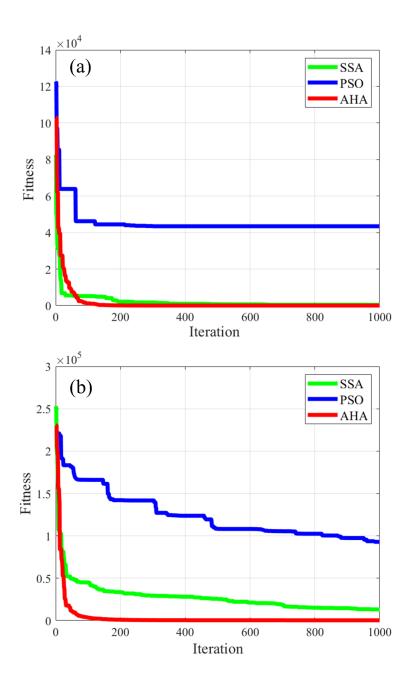




704



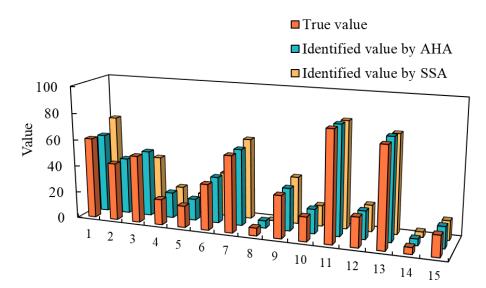




706 **Figure 7**







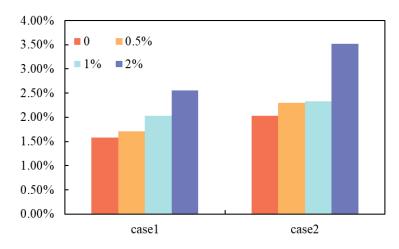
Variables to be identified

708

709 **Figure 8**







711 **Figure 9**





712 Table 1 Fundamental values and ranges of aquifer parameters.

Parameter	Value or range
Hydraulic conductivity of zone 1, K_1 (m/d)	(50,70)
Hydraulic conductivity of zone 2, K_2 (m/d)	(35,55)
Hydraulic conductivity of zone 3, K_3 (m/d)	(40,60)
Specific yield of zone 1, μ_1	0.27
Specific yield of zone 2, μ_2	0.22
Specific yield of zone 3, μ_3	0.25
Longitudinal dispersity of zone 1 (m)	40
Longitudinal dispersity of zone 2 (m)	30
Longitudinal dispersity of zone 3 (m)	35
Grid spacing in X and Y direction (m)	50
Recharge rate (m/d)	0.00042
Initial concentration (mg/L)	50
Length of the stress period (y)	10
Aquifer thickness(m)	10
Groundwater level at the western boundary, H_1 (m)	(18,20)
Groundwater level at the eastern boundary, $H_2(m)$	(15,17)

https://doi.org/10.5194/egusphere-2025-2083 Preprint. Discussion started: 20 May 2025 © Author(s) 2025. CC BY 4.0 License.





714

Table 2 Ranges of values of pollution sources.

Pollution		R	elease inten	sity (g/d×10	0)	
sources	T1	T2	Т3	T4	T5	T6-T10
S1	(0,86.4)	(0,69.12)	(0,60.48)	(0,51.84)	(0,43.2)	0
S2	(0,103.68)	(0,86.4)	(0,77.76)	(0,69.12)	(0,51.84)	0





716 Table 3 True values of the variables to be identified.

Variables to be identified	True value
K_1 (m/d)	60.37
K_2 (m/d)	42.84
K_3 (m/d)	50.17
H_1 (m)	19.09
$H_2(\mathbf{m})$	16.11
S_1T_1 (g/d×10)	34.25
$S_1T_2(g/d\times10)$	57.07
S_1T_3 (g/d×10)	57.99
$S_1T_4(g/d\times10)$	31.76
$S_1T_5(g/d\times10)$	18.14
$S_2T_1\left(g/d\times10\right)$	82.07
$S_2T_2\left(g/d\times10\right)$	22.18
S_2T_3 (g/d×10)	74.35
$S_2T_4(g/d\times10)$	49.24
$S_2T_5(g/d\times10)$	15.84





Table 4 Fundamental values and ranges of aquifer parameters and pollution

719 sources.

Parameter	Value or range
Specific yield	0.24
Transverse dispersity (m)	9.8
Longitudinal dispersity (m)	40
Aquifer thickness(m)	40
Grid spacing in x-direction(m)	20
Grid spacing in y-direction(m)	20
Number of stress periods	5
Hydraulic conductivity(m/d)	(30,50)
Fluxes of contamination source during stress period(g/d)	(0,52)

720





721 Table 5 True values of the variables to be identified.

Variables to be identified	True value
K_1 (m/d)	45.93
K_2 (m/d)	46.54
K_3 (m/d)	32.11
K_4 (m/d)	44.23
$S_1T_1(g/d)$	38.05
$S_1T_2(g/d)$	32.24
$S_1T_3(g/d)$	24.96
$S_1T_4(g/d)$	5.17
$S_1T_5(g/d)$	25.42
$S_2T_1(g/d)$	31.15
$S_2T_2(g/d)$	39.94
$S_2T_3(g/d)$	51.5
$S_2T_4(g/d)$	49.47
$S_2T_5(g/d)$	31.53
$S_3T_1(g/d)$	27.49
$S_3T_2(g/d)$	26.93
$S_3T_3(g/d)$	5.95
$S_3T_4(g/d)$	30.5
$S_3T_5(g/d)$	23.7





723 Table 6 A comparison of the accuracies of the assessed surrogate models.

Case	Surrogate model	\mathbb{R}^2	MRE	RMSE
C1	Kriging	0.9942	13.43%	11.8262
Case1	BPNN	0.9994	3.70%	3.6526
Case2	Kriging	0.9837	9.98%	37.7547
	BPNN	0.9989	4.48%	9.8488

724





Table 7 A comparison of inversion values under different noise levels for case

726 **study 1.**

Unknown	True	Inversion values under different noise levels							
variables	value	0	0.5%	1%	2%	Relative	eerror		
K_1	60.37	58.91	59.46	61.16	61.15	2.42%	1.50%	1.31%	1.29%
K_2	42.84	42.12	41.73	41.72	42.18	1.67%	2.58%	2.61%	1.54%
K_3	50.17	49.28	48.52	48.58	50.01	1.78%	3.29%	3.17%	0.31%
H_1	19.09	19.10	19.04	19.06	19.27	0.06%	0.24%	0.18%	0.96%
H_2	16.11	16.05	15.97	16.01	16.27	0.40%	0.87%	0.64%	0.97%
S_1T_1	34.25	34.65	34.82	35.37	36.50	1.16%	1.66%	3.26%	6.57%
S_1T_2	57.07	57.20	57.35	57.66	58.79	0.24%	0.49%	1.04%	3.01%
S_1T_3	5.80	5.48	5.59	5.64	5.56	5.49%	3.63%	2.78%	4.19%
S_1T_4	31.76	31.80	31.84	31.99	32.71	0.15%	0.25%	0.74%	3.00%
S_1T_5	18.14	18.21	18.24	18.31	18.63	0.39%	0.55%	0.96%	2.73%
S_2T_1	82.07	81.45	81.67	82.48	84.62	0.76%	0.50%	0.49%	3.10%
S_2T_2	22.18	21.02	20.99	21.10	21.86	5.22%	5.37%	4.87%	1.44%
S_2T_3	74.35	75.69	75.95	76.44	77.69	1.80%	2.15%	2.81%	4.49%
S_2T_4	4.92	4.86	4.85	4.74	4.84	1.37%	1.48%	3.76%	1.78%
S_2T_5	15.84	15.95	16.00	16.12	16.29	0.73%	1.06%	1.81%	2.86%





$\begin{tabular}{ll} \textbf{Table 8 A comparison of inversion values under different noise levels for case} \\ \end{tabular}$

729 **study 2.**

Unknow	Тті	Inversi	on value	s under		noise lev	els		
n variables	True value	0	0.5%	1%	2%	Relative	e error		
K_1	45.93	44.94	45.44	45.07	46.01	2.15%	1.07%	1.87%	0.17%
K_2	46.54	46.68	47.28	46.83	47.92	0.29%	1.59%	0.62%	2.97%
K_3	32.11	32.08	31.91	32.05	31.73	0.08%	0.62%	0.20%	1.19%
K_4	44.23	44.56	43.79	44.35	42.95	0.75%	0.98%	0.26%	2.89%
S_1T_1	38.05	37.48	37.59	37.85	38.14	1.48%	1.22%	0.51%	0.23%
S_1T_2	32.24	32.84	32.55	33.10	32.42	1.84%	0.95%	2.65%	0.55%
S_1T_3	24.96	26.75	26.46	26.89	26.48	7.18%	6.01%	7.74%	6.09%
S_1T_4	5.17	4.89	4.85	4.93	4.77	5.44%	6.33%	4.79%	7.82%
S_1T_5	25.42	26.48	26.29	26.69	26.42	4.18%	3.43%	5.03%	3.94%
S_2T_1	31.15	31.17	31.21	31.38	31.48	0.08%	0.19%	0.74%	1.07%
S_2T_2	39.94	40.17	40.12	40.65	40.58	0.57%	0.43%	1.76%	1.59%
S_2T_3	51.5	51.77	51.74	52.00	52.00	0.53%	0.47%	0.97%	0.97%
S_2T_4	49.47	48.91	48.81	49.51	49.36	1.13%	1.33%	0.09%	0.21%
S_2T_5	31.53	33.54	33.30	33.41	33.03	6.38%	5.61%	5.97%	4.75%
S_3T_1	27.49	27.61	28.03	28.01	28.75	0.43%	1.96%	1.90%	4.59%
S_3T_2	26.93	27.33	27.88	27.68	28.80	1.47%	3.52%	2.76%	6.95%
S_3T_3	5.95	5.97	6.14	6.11	6.38	0.27%	3.15%	2.66%	7.13%
S_3T_4	30.5	30.97	31.18	31.16	31.70	1.54%	2.21%	2.16%	3.92%
S_3T_5	23.7	23.05	24.32	24.06	26.06	2.77%	2.59%	1.49%	9.95%





731 Table 9 Mean relative errors of the two case studies under different noise levels.

	Differe	nt noise l	evels	
case	0	0.5%	1%	2%
case1	1.58%	1.71%	2.03%	2.55%
case2	2.03%	2.30%	2.33%	3.52%

732