

1 **Joint characterization of heterogeneous conductivity fields and pumping well**
2 **attributes through iterative ensemble smoother with a reduced-order modeling**
3 **strategy for solute transport**

4
5 Chuan-An Xia¹, Jiayun Li^{2*}, Bill X. Hu³, Alberto Guadagnini^{4,5}, Monica Riva^{4*}

6
7 ¹Zijin School of Geology and Mining, Fuzhou University, Fuzhou, China

8 ²Fujian Provincial Key Lab of Coastal Basin Environment, Fujian Polytechnic Normal
9 University, Fuqing, China

10 ³School of Water Conservancy & Environment, University of Jinan, Jinan, China

11 ⁴Dipartimento di Ingegneria Civile e Ambientale, Politecnico di Milano, Milano, Italy

12 ⁵Sonny Astani Department of Civil and Environmental Engineering, Viterbi School of
13 Engineering, Los Angeles, California 90089-2531, USA

14
15
16
17 Submitted to: *Hydrology and Earth System Sciences*

18
19 Corresponding author: Jiayun Li; Monica Riva

20 Email: lijy@fpnu.edu.cn; monica.riva@polimi.it

21

22

Abstract

23 We develop and test an efficient and accurate theoretical and computational
24 framework to jointly estimate spatially variable hydraulic conductivity and identify
25 unknown pumping well locations and rates in a two-dimensional confined aquifer.
26 The approach (denoted as iES_ROM) integrates an iterative Ensemble Smoother (iES)
27 with a Reduced-Order Model (ROM) for solute transport taking place across an
28 otherwise steady-state groundwater flow field. This offers a computationally efficient
29 alternative to the Full System Model (iES_FSM) upon addressing the high
30 computational demands of ensemble-based data assimilation methods, which typically
31 require large ensemble sizes to characterize uncertainties in (randomly) heterogeneous
32 aquifers. Our iES_ROM is constructed through proper orthogonal decomposition. It is
33 then evaluated across a collection of 28 test cases exploring variations in model
34 dimension, ensemble size, measurement noise, monitoring network, and statistical
35 properties of the (underlying randomly heterogeneous) conductivity field. Our results
36 support the ability of iES_ROM to accurately estimate conductivity and identify
37 pumping well attributes under diverse configurations, attaining a quality of
38 performance similar to iES_FSM. When using moderate ROM dimensions ($n = 25-30$)
39 and ensemble size (i.e., 500-1000), the accuracy of iES_ROM does not vary
40 significantly while computational time is reduced by nearly an order of magnitude.
41 Our approach thus provides a reliable and cost-effective tool for inverse modeling in
42 groundwater systems with uncertain parameters.

43 Keywords: reduced-order model; proper orthogonal decomposition; iterative

44 ensemble smoother; pumping well identification; groundwater

45 **1. Introduction**

46 Assessment of groundwater flow and transport scenarios is typically plagued by
47 uncertainties associated with model structure and parametrization. A major source of
48 uncertainty often examined concerns the poorly constrained assessment of pollution
49 sources. Our ability to identify spatial locations of these sources exerts significant
50 influence on the design of contaminant monitoring, management, and remediation
51 strategies. Contaminant release to an aquifer is characterized through the spatial
52 location of sources, the temporal variability of release fluxes, and solute
53 concentrations involved (Chen et al., 2018; Xu et al., 2018; Mo et al., 2019).
54 Uncertainties linked to groundwater abstraction scheduling also play a critical role, as
55 operational details of pumping wells are not always fully documented. For example,
56 this might correspond to a scenario where such information is not disclosed to ensure
57 privacy protection or uncertainties are induced by geocoding practices and/or
58 measurement devices. Further to this, in some regions groundwater may be accessed
59 through wells that are not officially registered or fully documented by industrial
60 operators and/or local residents. Despite the relevance of these issues, only limited
61 research has been devoted to the identification and quantification of pumping rates
62 and spatial locations of such hidden wells.~~Despite the relevance of these issues, only~~
63 ~~limited research has been devoted to the identification and quantification of pumping~~
64 ~~rates and spatial locations of such hidden wells.~~

65 In this broad context, we recall that a considerable body of research has focused
66 on estimating key parameters (such as hydraulic conductivity) in groundwater flow

67 and transport models through ensemble-based Data Assimilation (DA) techniques
68 (Chen and Zhang, 2006; Tong et al., 2013; Chen and Oliver, 2013; Zhang et al., 2018;
69 Xia et al., 2018, 2024). These approaches aim at enhancing the accuracy of simulated
70 system states (e.g., hydraulic heads and solute concentrations). While their capability
71 to jointly estimate parameters and update system states has been broadly explored,
72 their high computational cost still constitutes a persistent limitation to their practical
73 routine use. This challenge primarily stems from the requirement for a large number
74 of realizations to ensure statistical convergence of Monte Carlo (MC) simulations
75 (e.g., Ballio and Guadagnini, 2004) in the forecast step of the DA process, and to
76 achieve reliable parameter estimates in the analysis step. The computational burden
77 becomes particularly significant when the selected model describing the system
78 behavior (hereafter termed as Full System Model (FSM)) must be repeatedly executed
79 for systems characterized by strong nonlinearities or requiring high (space-time)
80 resolution of state variables and parameters.

81 To alleviate these computational constraints, recent studies explore the benefit of
82 relying on surrogate (or reduced-order) models that approximate the behavior of the
83 full system while maintaining sufficient accuracy for inverse modeling workflows and
84 Uncertainty Quantification (UQ).

85 In this framework, efforts to mitigate computational limitations of
86 (ensemble-based) DA methods primarily focus on the adoption of localization
87 techniques (e.g., Xia et al., 2018, 2024; Luo and Bahkta, 2020) or surrogate modeling
88 strategies (e.g., Zhang et al., 2018; Mo et al., 2019 and references therein). Main

89 advantages associated with location approaches are related to the observation that
90 they (i) substantially reduce computational costs upon requiring only a limited
91 number of Monte Carlo realizations of the FSM, while maintaining acceptable
92 accuracy of the assimilated results, and (ii) retain a physically-based and
93 mathematically-tractable formulation. As a notable drawback of these approaches, we
94 note that the value of the information associated with diverse measurements may be
95 partially suppressed due to the use of distance- or correlation-based localization,
96 which might constrain the strength of the spatial influence of observations. As a
97 consequence, the ensuing (empirical/sample) probability density functions (PDFs) of
98 model parameters and system states often display reduced accuracy and fail to fully
99 capture the underlying uncertainty structure. To mitigate these limitations, an
100 alternative line of research explores the use of surrogate models (SMs), which aim at
101 emulating the response of the Full System Model with significantly reduced
102 computational cost while preserving the salient physics of the system.

103 Surrogate models are rapidly emerging as a promising complement to FSMs for
104 reducing computational burdens associated with the forecast steps of ensemble-based
105 DA procedures. Among the various SM strategies, data-driven approaches based on
106 machine learning (e.g., Ju et al., 2018) and deep learning (e.g., Mo et al., 2019) can be
107 employed for emulating groundwater flow and transport processes taking place in
108 heterogeneous media. For example, Ju et al. (2018) rely on Gaussian Process
109 regression to describe relationships between the coefficients of a Karhunen-Loève
110 (KL) expansion (employed to characterize a spatially heterogeneous hydraulic

111 conductivity field) and (point-wise) simulated observations. This approach is shown
112 to achieve approximately an order of magnitude reduction in computational time as
113 compared with the standard iterative Ensemble Smoother (iES). Otherwise, this gain
114 in efficiency is associated with a reduced accuracy in simulated hydraulic heads,
115 which in turn compromises the reliability of the estimated conductivity field. Mo et al.
116 (2019) employ deep autoregressive neural networks as an FSM surrogate to
117 reconstruct conductivity fields and identify contaminant source characteristics.
118 However, their approach still requires a significant computational effort, as it heavily
119 relies on a high number (about 1,500 in their exemplary setting) of MC realizations of
120 the FSM for network training. While these studies show a clear potential of
121 data-driven surrogates for accelerating DA workflows, they also highlight the need for
122 a fundamental trade-off between computational efficiency and model accuracy, thus
123 underscoring the potential value of alternative surrogate modeling strategies.

124 In contrast to data-driven models, that typically operate as black-box
125 representations, projection-based Reduced-Order Models (ROMs) are physics-based
126 (e.g., Razavi et al., 2012; Asher et al., 2015; Chen et al., 2017; Xia et al., 2020, 2025).
127 ROMs are typically constructed upon projecting the governing equations and
128 boundary conditions of the onto a lower-dimensional subspace spanned by a set of
129 basis functions. The latter are commonly derived through, e.g., Proper Orthogonal
130 Decomposition (POD) of multiple FSM solutions, referred to as *snapshots*. This
131 procedure effectively reduces the dimensionality of the system state space. The
132 random field representing the system state can then be expressed as a linear

133 combination of the dominant eigenfunctions obtained from the Fredholm integral
134 equation associated with the covariance matrix of the snapshots. Leading
135 eigenfunctions are then identified as the basis functions defining the reduced subspace.
136 Substantial computational savings are then achieved upon resting on the solution of
137 the ensuing low-dimensional linear system. When implemented in the context of
138 numerical MC frameworks, the collection of ROM-generated solutions constitutes
139 what is commonly referred to as a Reduced-Order Monte Carlo (ROMC) simulation
140 framework.

141 Reduced-order modeling has received growing attention in the context of
142 groundwater flow (Pasetto et al., 2011, 2013, 2014; Li et al., 2013a; Boyce et al., 2015;
143 Stanko et al., 2016; Xia et al., 2020, 2025) and solute transport (Luo et al., 2012; Li et
144 al., 2013b; Rizzo et al., 2018) scenarios. Its potential is evidenced across a wide range
145 of hydrogeological configurations, including confined (e.g., Pasetto et al., 2011) and
146 unconfined (e.g., Stanko et al., 2016) aquifer systems, homogeneous (e.g., Li et al.,
147 2013a) and heterogeneous (e.g., Pasetto et al., 2013) media, as well as scenarios with
148 (e.g., Xia et al., 2020) or without (e.g., Pasetto et al., 2014) pumping wells operating
149 therein. Several studies further advance development of ROMC strategies for UQ in
150 groundwater flow modeling. Pasetto et al. (2014) show that the accuracy of UQ
151 results relying on ROMC in the presence of steady-state groundwater flow strongly
152 depends on the quality and the number of snapshots, the latter directly influencing
153 representativeness of the basis functions. To mitigate this limitation, Xia et al. (2020)
154 propose deriving basis functions as the leading eigenvectors of (second-order)

155 approximations of hydraulic head covariances. The latter are obtained upon solving
156 the associated moment equations for steady-state groundwater flow (Zhang and Lu,
157 2002; Xia et al., 2019). Even as reduction of the dimensionality of the head space
158 provides substantial computational savings, projection of the basis functions onto the
159 ensuing (typically large) system matrix remains computationally intensive, thereby
160 still constituting a limiting factor to efficiency gains. Xia et al. (2025) address this
161 challenge by extending their approach to perform dimensionality reduction for both
162 (spatially variable) transmissivity and hydraulic head fields in a steady-state
163 groundwater flow setting and achieving additional computational savings while
164 maintaining high accuracy. Despite these advancements, most existing ROM and
165 ROMC approaches are still fraught with difficulties in efficiently capturing strongly
166 nonlinear system dynamics and adapting to evolving state conditions, underscoring
167 the need for more flexible and computationally efficient reduced-order frameworks.

168 With reference to solute transport, ROMs have been developed for both
169 homogeneous (Luo et al., 2012) and heterogeneous (Li et al., 2013b; Rizzo et al.,
170 2018) aquifer systems. Li et al. (2013a) further consider construction of ROMs to
171 tackle density-dependent groundwater flow taking place across homogeneous and
172 heterogeneous domains. Otherwise, studies explicitly focusing on the development of
173 ROMC approaches for UQ of solute transport remain limited. [Although conceptual
174 insights can be drawn from ROMC studies addressing groundwater flow \(e.g., Pasetto
175 et al., 2014; Xia et al., 2020, 2025\), influence of key factors \(such as, e.g.,
176 dimensionality of the reduced concentration space and strength of hydraulic](#)

177 conductivity heterogeneity) on accuracy and robustness of ROMC-based UQ still
178 remains poorly characterized. Although conceptual insights can be drawn from ROMC
179 studies addressing groundwater flow (e.g., Pasetto et al., 2014; Xia et al., 2020, 2025),
180 influence of key factors (such as, e.g., dimensionality of the reduced concentration
181 space and strength of hydraulic conductivity heterogeneity) on accuracy and
182 robustness of ROMC-based UQ still remains poorly characterized.

183 Building upon these works, the present study introduces a novel framework that
184 integrates the iES with a ROM for solute transport (hereafter referred to as
185 iES_ROM). Building upon these works, the present study introduces a novel
186 framework that integrates the iES with a ROM for solute transport (hereafter referred
187 to as iES_ROM). The ensuing framework enables one to efficiently quantify
188 uncertainty and jointly estimate system parameters in groundwater-related modeling
189 scenarios. The proposed method is then applied to simultaneously identify pumping
190 rate and spatial location of (otherwise hidden) wells operating within the system,
191 while providing estimates of the spatially heterogeneous hydraulic conductivity field
192 under conditions of steady-state flow and transient solute transport. In the iES_ROM
193 framework, the steady-state flow field is evaluated through the FSM, whereas the
194 transient solute transport is represented by a computationally efficient ROM. The
195 required snapshots and associated POD are generated only once. These are
196 subsequently employed throughout the entire DA process, thus avoiding repeated
197 high-fidelity simulations. To ensure transparent benchmarking, the performance of
198 iES_ROM is systematically compared with that of a reference approach (termed

199 iES_FSM) which relies entirely on the FSM associated with synthetic scenarios.
200 Comparative analyses are performed across a variety of synthetic scenarios,
201 encompassing diverse ROM dimensions, ensemble sizes, measurement qualities and
202 quantities, as well as distinct statistical descriptors of the initial conductivity ensemble
203 and snapshot sizes.

204 The study is organized as follows. Section 2 introduces the theoretical
205 background of groundwater flow and solute transport and details the integration of
206 ROMC simulation within the iES framework. Section 3 describes the test cases
207 designed to evaluate the proposed approach. Section 4 illustrates and discusses the
208 main results, and Section 5 summarizes the key findings.

209 **2. Theory background and methodology**

210 **2.1 Groundwater flow and solute transport**

211 We consider two-dimensional steady-state groundwater flow governed by:

$$212 \nabla \cdot [K(\mathbf{x}) \nabla h(\mathbf{x})] + q_s(\mathbf{x}) = 0 \quad (1)$$

213 where $\mathbf{x} = [x_1, x_2]$ is a vector of spatial coordinates in domain Ω^2 ; h is hydraulic
214 head; K is (isotropic) hydraulic conductivity; and q_s is a source/sink term. We
215 conceptualize K as a spatially heterogeneous random field, associated with a given
216 spatial correlation structure. The source/sink term in Equation (1) corresponds to a
217 production well associated with an uncertain pumping rate and location in the domain.
218 Propagation of uncertainty related to model parameters and/or forcing terms onto
219 hydraulic heads and fluxes is typically assessed through numerical Monte Carlo (MC)
220 simulations (see, e.g., Ballio and Guadagnini, 2004; Xia et al., 2020, 2024, and

221 references therein).

222 We consider (non-reactive) solute transport evolving in Ω^2 to be described
223 through:

$$224 \quad \nabla \cdot [D \nabla c(\mathbf{x}, t)] - \nabla(\mathbf{q}(\mathbf{x})c(\mathbf{x}, t)) + \frac{q_s(\mathbf{x})}{\theta} c_s(\mathbf{x}, t) = \frac{\partial c(\mathbf{x}, t)}{\partial t} \quad (2)$$

225 Here, t denotes time; c is solute concentration; D is the (isotropic) dispersion
226 coefficient; θ is effective porosity; c_s is solute concentration corresponding to q_s ;
227 and $\mathbf{q}(\mathbf{x}) = -(K(\mathbf{x})/\theta)\nabla h(\mathbf{x})$ is an effective velocity associated with solute
228 transport.

229 Numerical methods (e.g., finite differences or finite elements) are commonly
230 employed to discretize Equations (1) and (2) that are then solved within a numerical
231 MC context. The probability distribution of state variables of interest (e.g., heads or
232 concentrations) is then evaluated at N nodes of an aptly designed numerical grid.
233 Consistent with Section 1, we refer to the model corresponding to the numerical
234 solution of the above equations as the Full System Model (FSM). When the domain is
235 characterized by a large spatial extent and/or one is interested in exploring the system
236 behavior across long temporal windows, performing numerical MC simulations
237 relying on FSM is associated with a heavy computational burden. To circumvent this
238 issue, we rely on the development and implementation of a Reduced-Order Model
239 (ROM) strategy for solute transport. We note that in this study we employ ROM
240 solely for solute transport because only limited computational costs are associated
241 with the steady-state flow condition we consider, as opposed to simulating transport.
242 Hereafter, we refer to numerical MC analyses grounded on ROM as ROMC.

243 2.2 Numerical Monte Carlo simulation framework for solute transport

244 2.2.1 Monte Carlo simulation setting for the Full System Model

245 We rely on a standard finite element method to solve the FSM described in
246 Section 2.1. When considering a total simulation time T_s , we express the linear
247 system associated with the numerical solution of solute transport through FSM within
248 time interval $[t, t + \Delta t]$ as:

$$249 \mathbf{A}^i \mathbf{c}^i = \mathbf{F}^i \quad (3)$$

250 Here, superscript i refers to the i^{th} MC realization ($i = 1, \dots, N_{MC}$, N_{MC} being the
251 total number of MC simulations) of FSM; \mathbf{A} is the full-system stiffness matrix (of
252 size $N \times N$); \mathbf{c} is the vector (of size $N \times 1$) of solute concentration values; and \mathbf{F} is the
253 stress vector (of size $N \times 1$) whose entries encompass source/sink terms and initial and
254 boundary conditions.

255 2.2.2 Reduced-order Monte Carlo simulation framework

256 We construct a reduced-order model for solute transport by approximating the
257 solution of solute concentration for the i^{th} MC realization of FSM. Consistent with the
258 work of Xue and Xie (2007) and Pinnau (2008), one can approximate \mathbf{c}^i as:

$$259 \mathbf{c}^i \approx \sum_{j=1}^n \alpha_j^i \mathbf{p}_j = \mathbf{P} \boldsymbol{\alpha}^i \quad (4)$$

260 Here, $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n]$ is a matrix (of size $N \times n$, n being the dimension of the ROM)
261 collecting the n nodal basis functions that are here obtained through a Proper
262 Orthogonal Decomposition (POD) approach (see below); $\boldsymbol{\alpha}^i = [\alpha_1^i, \alpha_2^i, \dots, \alpha_n^i]^T$ (T
263 representing transpose) is a vector (of size $n \times 1$) of Fourier coefficients (Pinnau, 2008).
264 Note that Equation (4) is different from a typical Karhunen-Loève expansion of \mathbf{c}^i

265 (i.e., $\mathbf{c}^i \approx \langle \mathbf{c} \rangle + \sum_{j=1}^n \alpha_j^i \mathbf{p}_j = \langle \mathbf{c} \rangle + \mathbf{P} \boldsymbol{\alpha}^i$, see Equation (11) in Li et al., 2013b). As we
 266 illustrate in Section 2.3, relying on Equation (4) enables straightforward (i) coding
 267 and (ii) compatibility with the iterative Ensemble Smoother (IES).

268 Substituting Equation (4) into Equation (3) and imposing the residual of the
 269 model equation associated with the approximated solution to be orthogonal to the
 270 projection space defined through \mathbf{P} yields:

271
$$\mathbf{P}^T \mathbf{A}^i \mathbf{P} \boldsymbol{\alpha}^i \approx \mathbf{P}^T \mathbf{F}^i \quad (5)$$

272 Solving Equation (5) (which is a linear system of size n) yields $\boldsymbol{\alpha}^i$ for the i^{th}
 273 MC realization of our ROMC strategy. Note that, when $n \ll N$, the computational
 274 effort required by our ROMC is much less than that of the standard MC.

275 The basis functions forming the entries of \mathbf{P} are computed as the leading
 276 eigenvectors (corresponding to the highest eigenvalues) of the covariance of solute
 277 concentration evaluated through N_{sn} numerical solutions (i.e., $\mathbf{c}^1, \mathbf{c}^2, \dots$, and
 278 $\mathbf{c}^{N_{sn}}$) of the FSM. Here, $N_{sn} = m \times N_t$, where m is the number of MC realizations
 279 of hydraulic conductivity that are randomly sampled from the initial ensemble of Y
 280 fields, each yielding $N_t = T_s / \Delta t$ (Δt corresponding to a uniform time step)
 281 numerical solutions of Equation (2). The basis functions forming the entries of \mathbf{P} are
 282 computed as from the leading eigenvectors (corresponding to the highest eigenvalues)
 283 of the covariance of totally N_{sn} numerical solutions (i.e., $\mathbf{c}^1, \mathbf{c}^2, \dots$, and $\mathbf{c}^{N_{sn}}$)
 284 through FSM. We point out that $N_{sn} = m \times N_t$, where m is the number of MC
 285 realizations which are arbitrarily chosen, each yielding $N_t = T_s / \Delta t$ numerical
 286 solutions of Equation (2), with uniform length of time step Δt . The leading

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

287 eigenvectors are computed through the Singular Value Decomposition (SVD)
 288 approach, i.e.:

$$289 \quad \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \text{svd}(\mathbf{E}\mathbf{E}^T) \quad (6)$$

290 where $\mathbf{E} = 1/\sqrt{N_{sm}}[\mathbf{c}^1, \mathbf{c}^2, \dots, \mathbf{c}^{N_{sm}}]$; \mathbf{U} (of size $N \times N$) is the left singular matrix
 291 whose j^{th} column is the j^{th} eigenvector of matrix $\mathbf{E}\mathbf{E}^T$ corresponding to the j^{th} singular
 292 value, λ_j ; and $\mathbf{\Lambda} = \text{diag}([\lambda_1, \lambda_2, \dots, \lambda_N])$ whose entries are ranked in descending
 293 order.

294 2.3 Iterative ensemble smoother

295 We denote by $\mathbf{m} = [Y_1, Y_2, \dots, Y_N, \ln q_s, x_{1,q_s}, x_{2,q_s}]^T$ the vector (of size $P = N+3$)
 296 whose entries correspond to the uncertain model parameters (i.e., the log-conductivity,
 297 $Y = \ln K$, field) and flow rate and location of a pumping well. In case the pumping rate
 298 and location are known, then $\mathbf{m} = [Y_1, Y_2, \dots, Y_N]^T$ and $P = N$. We further denote by
 299 $\mathbf{d} = [d_1, d_2, \dots, d_o]^T$ the vector (of size O) of the observations (i.e., measured head
 300 and concentration values). To estimate \mathbf{m} , we implement the iES (Luo and Bhakta,
 301 2020; Xia et al., 2024):

$$302 \quad \begin{cases} \mathbf{m}^{k+1} = \mathbf{m}^k + \underline{\underline{\mathbf{K}}}_{\text{Gain}}^k \Delta \mathbf{d}^k \\ \underline{\underline{\mathbf{K}}}_{\text{Gain}}^k = \underline{\underline{\mathbf{S}}}_m^k (\underline{\underline{\mathbf{S}}}_d^k)^T (\underline{\underline{\mathbf{S}}}_d^k (\underline{\underline{\mathbf{S}}}_d^k)^T + \gamma^k \mathbf{I})^{-1} \quad \text{with } \gamma^i = \xi^i \text{trace}(\underline{\underline{\mathbf{S}}}_d^i (\underline{\underline{\mathbf{S}}}_d^i)^T) / O. \\ \Delta \mathbf{d}^k = g(\mathbf{m}^k) - \mathbf{d} \end{cases} \quad (7)$$

303 Here, superscript k is the index of the iteration step; matrices
 304 $\underline{\underline{\mathbf{S}}}_m^k = [\mathbf{m}_1^k - \bar{\mathbf{m}}^k, \dots, \mathbf{m}_N^k - \bar{\mathbf{m}}^k] / \sqrt{N-1}$ (of size $P \times N$, where $\bar{\mathbf{m}}^k = \sum_{j=1}^N \mathbf{m}_j^k / N$) and
 305 $\underline{\underline{\mathbf{S}}}_d^k = [g(\mathbf{m}_1^k) - g(\bar{\mathbf{m}}^k), \dots, g(\mathbf{m}_N^k) - g(\bar{\mathbf{m}}^k)] / \sqrt{N-1}$ (of size $O \times N$, where $g(\cdot)$
 306 represents model operator being either FSM or ROM) collect the ensemble anomalies

307 of parameters and simulated observations associated with the k^{th} iteration step; \mathbf{I} is
 308 the identity matrix (of size $O \times O$); and ξ^k is an adaptive coefficient (Luo et al., 2015)
 309 associated with each iteration of the Levenberg-Marquardt (LM; Levenberg, 1944)
 310 algorithm. We set $\xi^0 = 10$ in our showcase application examples (see Section 3) and
 311 follow Luo and Bhakta (2020) to update its value for the remaining iteration steps.

312 In the case of $g(\cdot)$ representing the model operator of ROM, we note that the
 313 approximation of solute concentration relying on Equation (4) is compatible with the
 314 implementation of Equation (7). The degree of compatibility of ROM to iES is
 315 reduced when considering a typical Karhunen-Loève expansion of \mathbf{c}^i (i.e.,
 316 $\mathbf{c}^i \approx \langle \mathbf{c} \rangle + \sum_{j=1}^n \alpha_j^i \mathbf{p}_j = \langle \mathbf{c} \rangle + \mathbf{P} \boldsymbol{\alpha}^i$). This is related to the observation that $\langle \mathbf{c} \rangle$ evolves
 317 with time and needs to be evaluated at each time step. This, in turn, implies that m
 318 numerical solutions of solute concentration through FSM need to be obtained to
 319 evaluate $\langle \mathbf{c} \rangle$ at every outer iteration of iES. Hence, computational advantages of
 320 employing ROM are reduced while coding complexity increases.

~~The degree of compatibility of ROM to iES is reduced when considering a typical Karhunen-Loève expansion of \mathbf{c}^i (i.e., $\mathbf{c}^i \approx \langle \mathbf{c} \rangle + \sum_{j=1}^n \alpha_j^i \mathbf{p}_j = \langle \mathbf{c} \rangle + \mathbf{P} \boldsymbol{\alpha}^i$). This is related to the observation that $\langle \mathbf{c} \rangle$ evolves with time and needs to be evaluated at each time step. This, in turn, implies that m numerical solutions of solute concentration through FSM need to be obtained to evaluate $\langle \mathbf{c} \rangle$ at every outer iteration of iES. Hence, computational advantages of employing ROM are reduced while coding complexity increases.~~ When approximating solute concentration via Equation (4), we only obtain
 327 m numerical solutions of solute concentration through FSM at the first outer iteration
 328

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

Field Code Changed

Formatted: Font: 12 pt, Font color: Auto

329 of iES. Leading eigenvectors are computed upon relying on these solutions and are
 330 then stored. The Fourier coefficients \mathbf{a}^i associated with time interval $[t, t + \Delta t]$ for
 331 each MC realization starting from the second outer iteration of iES are obtained solely
 332 through solving Equation (5).

333 When implementing the LM algorithm during optimization, we set both the inner
 334 and outer iteration numbers equal to 10 (see also Luo and Bhakta, 2020). Additionally,
 335 a stopping criterion $(\delta_{k-1} - \delta_k) / \delta_{k-1} \times 100\% \leq 10^{-6}$ (where
 336 $\delta_k = \frac{1}{N} \sum_{j=1}^N \left\{ \left(\mathbf{d}_j^k - g(\mathbf{m}_j^k) \right)^T \mathbf{C}_d^{-1} \left(\mathbf{d}_j^k - g(\mathbf{m}_j^k) \right) \right\}$), is set.

337 **2.4 Implementation and Computational cost**

338 We denote by iES_FSM and iES_ROM the approaches associated with coupling
 339 the iES with FSM and ROM, respectively. [A workflow for iES ROM is summarized](#)
 340 [in Fig. 1](#). The total number of MC realizations is denoted as N_{MC} . Neglecting the
 341 computational cost of the inner iterations and assuming iES comprises N_{out} outer
 342 iterations, the main computational costs of either method can be divided into two
 343 components, corresponding to forecast and analysis step (Table 1), respectively. In the
 344 forecast step, a number of $(N_{out} + 1)$ MC simulations for groundwater flow and solute
 345 transport are required. Otherwise, Equation (7) is evaluated N_{out} times in the
 346 analysis step. The steady-state groundwater flow is solved through the FSM in both
 347 iES_FSM and iES_ROM, with a computational cost of order $O(N^3 N_{MC})$. The main
 348 computational cost for the N_{MC} FSM-based MC realizations of solute transport at a
 349 single time step in iES_FSM is $O(N^3 N_{MC})$, while being $O((sN + N^2) N_{MC})$
 350 (where $s \approx 7$ or ≈ 15 in two and three dimensions, respectively) for iES_ROM. These

351 computational efforts correspond to the projection of the full-system stiffness matrix
352 onto the reduced-order space of the system state (i.e., solute concentration).
353 Computational costs associated with solving Equation (7) coincide for both
354 approaches and are here denoted as C_8 . We further note that, with reference to
355 iES_ROM, the N_{sn} solutions of solute concentration obtained through FSM
356 (associated with a computational cost of order $O(N^3 N_{sn})$) and the basis functions
357 obtained through SVD (with a computational cost of order $O(n N_{sn}^2)$) are calculated
358 only once and stored. When the grid mesh employed is large or the simulation time is
359 long, computational savings through iES_ROM compared with iES_FSM become
360 significant.

361 **3. Exemplary scenarios**

362 We consider a two-dimensional computational domain of size 4×2 to simulate a
363 synthetic sandbox-scale experiment where (non-reactive) solute transport under
364 steady-state flow is considered (see Fig. 42). Here and hereafter, all quantities are
365 given in consistent (length/mass/time) units. Concerning groundwater flow, the left
366 and right sides of the domain are associated with constant head boundary conditions
367 with $H = 3$ and 2, respectively. The top and bottom sides correspond to boundary
368 conditions of no flow. A pumping well with an unknown pumping rate and location is
369 considered in the setting. A fixed concentration boundary is set at point $(0, 1)$ (see red
370 triangle in Fig. 42) with a constant concentration of 100, while the initial
371 concentration across the domain is set to zero. We use the standard finite element
372 method to obtain the numerical solutions of head and concentration. The numerical

373 mesh adopted comprises $41 \times 21 = 861$ nodes and 1,600 triangle elements. [A uniform](#)
 374 [time step of 1 day is considered, our analyses encompassing a total simulation time of](#)
 375 [10 days \(i.e., \$T_s = 10\$ days and \$N_t = 10\$ \).](#) ~~A uniform time step of 1 day is considered for~~
 376 ~~a total simulation time of 10 days.~~

377 The logarithm of conductivity ($Y = \ln K$) is considered as a spatially
 378 heterogeneous (correlated) Gaussian random field with an exponential covariance
 379 function (C_Y) given by:

$$380 \quad C_Y = \sigma_Y^2 \exp\left(-\left(\frac{d_{x_1}}{\lambda_{x_1}} + \frac{d_{x_2}}{\lambda_{x_2}}\right)\right) \quad (8)$$

381 where σ_Y^2 is the variance of Y ; d_{x_i} ($i = x, y$) is separation (lag) distance between
 382 two given points in the i -direction; λ_{x_i} (with $i = x, y$) is the correlation length of Y in
 383 the i -direction. The corresponding mean of Y is denoted as μ . The initial ensemble of
 384 Y fields is synthetically generated through the well-known and widely tested GSLIB
 385 suite (Deutsch and Journel, 1998) upon setting λ_{x_1} and λ_{x_2} equal to 1.0 and 0.5,
 386 respectively. The reference Y field (Fig. 1a) is generated upon setting $\mu = 0.8$, σ_Y^2
 387 $= 1.0$, $\lambda_{x_1} = 1.0$, and $\lambda_{x_2} = 0.5$.

388 The pumping rate (i.e., q_s), x_1 and x_2 -coordinates (denoted as x_{1,q_s} and
 389 x_{2,q_s} , respectively) of the pumping location are considered to be random variables,
 390 each associated by a Gaussian distribution. The gray zone in Fig. [4b-2b](#) encompasses
 391 the possible locations where a pumping well is operating. The initial collection
 392 (ensemble) of values of q_s , x_{1,q_s} , and x_{2,q_s} and their reference counterparts are
 393 sampled from Gaussian distributions characterized by mean (standard deviation)
 394 equal to 0.50 (0.25), 1.00 (0.25), and 1.00 (0.25), respectively. These settings ensure

395 that the randomly generated samples of x_{1,q_s} and x_{2,q_s} are mostly within the
396 coordinate ranges indicated by the gray zone in Fig. 2b. Reference values are $q_s =$
397 1.03, $x_{1,q_s} = 1.38$, and $x_{2,q_s} = 1.40$ (see Fig. 4b2b, red cross symbol). Figure 4e-2c
398 depicts the simulated head field associated with the reference conductivity field,
399 pumping rate, and location. Figure 1d depicts simulated concentrations at the final
400 simulation time. Observations, including (steady-state) head and solute concentration
401 at each time step, are collected at a number (denoted as N_m) of monitoring wells
402 distributed across the aquifer according to some pre-defined patterns (Fig. 4b2b-d).
403 Each measurement is taken as the sum of the simulated head (or concentration) and a
404 white noise with zero mean and standard deviation equal to σ_{obs} .

405 To explore the potential of iES_ROM, several showcases are designed to
406 highlight key features of interest. Five groups of test cases (TCs) are designed and
407 organized as detailed in the following (see also Table 1).

- 408 ➤ **Group A.** It includes twelve TCs (i.e., TC1-TC12), enabling us to compare
409 performances of iES_FSM and iES_ROM associated with diverse values of
410 n when the pumping rate and locations are either known (TC1-TC6) or
411 unknown (TC7-TC12). The dimension of the ROM is considered equal to $\{5,$
412 $10, 15, 20, 25, 30\}$, these values being consistent with those most commonly
413 analyzed in previous studies (Pasetto et al., 2014; Xia et al., 2020, 2025).
- 414 ➤ **Group B.** It includes four TCs (i.e., TC6 and TC13-TC15), enabling us to
415 compare the performances of iES_FSM and iES_ROM with the largest
416 value of n analyzed (i.e., $n = 30$) and considering diverse values of N_{MC}

417 corresponding to {30, 100, 500, 10,000}. The latter are values of N_{MC}
418 commonly tested in previous studies (Chen and Zhang, 2006; Xia et al.,
419 2021, 2024).

420 ➤ **Group C.** It includes five TCs (i.e., TC6 and TC16-TC19), designed to
421 analyze the ability of iES_ROM to cope with diverse quality and quantity of
422 available measurements. Performances of iES_FSM and iES_ROM are also
423 compared when $\sigma_{obs} = \{0.001, 0.01, 0.1\}$ and the number of observation
424 locations corresponds to a value selected from {9 (Fig. 4e2b), 18 (Fig. 4e2c),
425 55 (Fig. 4e2d)}.

426 ➤ **Group D.** It includes five TCs (i.e., TC6 and TC_20-TC23), enabling us to
427 study the effect of μ and σ_Y^2 of the initial ensemble of Y on the
428 accuracies of estimates of conductivity and pumping rate and well location
429 through iES_FSM and iES_ROM. Values of μ and σ_Y^2 of the initial
430 ensemble of Y fields are selected from {-0.5, 1.2, 2.0} and {0.01, 1.0, 2.0},
431 respectively.

432 ➤ **Group E.** It includes six TCs (i.e., TC6 and TC24-TC28), with the aim of
433 investigating the effect of N_{sn} on the accuracies of the estimation of
434 conductivity and well pumping rate and location through iES_ROM and on
435 computation time requirements. Values of N_{sn} in TC24-TC28 and TC6 are
436 equal to 30, 100, 300, 500, 1,000, and 10,000, respectively.

437 Note that, without specified otherwise, default settings for the above mentioned
438 TCs correspond to TC6 which is designed with $n = 30$, $N_{MC} = 10,000$, $N_{sn} = 10,000$,

439 $N_m = 55$, $\sigma_{obs} = 0.01$, and values of μ and σ_Y^2 of the initial ensemble of Y equal
 440 to 1.2 and 1.0, respectively. Except for TC8-TC12, the source/sink term is associated
 441 with uncertainty.

442 To quantify the accuracy of conductivity estimates through iES_ROM and
 443 iES_FSM, we consider absolute error between estimated and reference values of Y
 444 (denoted as E_Y) and estimate of the standard deviation (denoted as S_Y) which are
 445 defined as:

$$446 \quad E_Y = \frac{1}{N} \sum_{i=1}^N \left| \langle Y_i \rangle^{est} - Y_i^{ref} \right|; \quad S_Y = \sqrt{\frac{1}{N} \sum_{i=1}^N (\sigma_{Y,i}^2)^{est}} \quad (9)$$

447 where $\langle Y_i \rangle^{est}$, $(\sigma_{Y,i}^2)^{est}$, and Y_i^{ref} denote estimated (ensemble) mean and variance,
 448 and reference value of Y at the i^{th} cell of the numerical grid, respectively.

449 Absolute errors and estimates of the standard deviations of $\ln q_s$, x_{1,q_s} , and
 450 x_{2,q_s} are employed to quantify the accuracy of the estimate of the pumping rate and
 451 well location:

$$452 \quad E_{q_s} = \left| \langle \ln q_s \rangle^{est} - \ln q_s^{ref} \right|; \quad E_{x_1} = \left| \langle x_{1,q_s} \rangle^{est} - x_{1,q_s}^{ref} \right|; \quad E_{x_2} = \left| \langle x_{2,q_s} \rangle^{est} - x_{2,q_s}^{ref} \right| \quad (10)$$

453 where $\langle \ln q_s \rangle^{est}$, $\langle x_{1,q_s} \rangle^{est}$, and $\langle x_{2,q_s} \rangle^{est}$ indicate estimated (ensemble) mean values
 454 of $\ln q_s$, x_{1,q_s} , and x_{2,q_s} respectively; and q_s^{ref} , x_{1,q_s}^{ref} , and x_{2,q_s}^{ref} are the reference
 455 values of q_s , x_{1,q_s} , and x_{2,q_s} , respectively. Estimates of the standard deviations of
 456 $\ln q_s$, x_{1,q_s} , and x_{2,q_s} are:

$$457 \quad S_{q_s} = \sqrt{(\sigma_{\ln q_s}^2)^{est}}; \quad S_{x_1} = \sqrt{(\sigma_{x_{1,q_s}}^2)^{est}}; \quad S_{x_2} = \sqrt{(\sigma_{x_{2,q_s}}^2)^{est}} \quad (11)$$

458 where $(\sigma_{\ln q_s}^2)^{est}$, $(\sigma_{x_{1,q_s}}^2)^{est}$, and $(\sigma_{x_{2,q_s}}^2)^{est}$ denote estimated (ensemble) variances of
 459 $\ln q_s$, x_{1,q_s} , and x_{2,q_s} , respectively.

460 As an additional metric, we then rely on the average absolute difference between
461 available data and model results:

$$462 \quad E_{obs} = \frac{1}{O} \sum_{i=1}^O \left| \langle d_i \rangle^{up} - d_i^{ref} \right| \quad (12)$$

463 where $\langle d_i \rangle^{up}$ and d_i^{ref} correspond to the (updated) result of the simulation process
464 and its reference observed counterpart at the i^{th} sampled location, respectively.

465 **4. Results and discussion**

466 **4.1 Impact of the dimension of the reduced-order model (Group A)**

467 Figure 2-3 depicts E_Y (Fig. 2a3a), S_Y (Fig. 2b3b), and E_{obs} (Fig. 2e3c)
468 versus the number of outer iterations for test cases (TCs) 1-6 obtained through
469 iES_ROM and iES_FSM, when well pumping rate and location are uncertain. Note
470 that results obtained through iES_FSM are independent of n (and are identical among
471 TCs 1-6) and are taken as references. Percentage differences (denoted as ΔE_Y)
472 between the values of E_Y obtained through iES_ROM and iES_FSM are depicted in
473 Fig. 2d3d. Corresponding results associated with percentage differences between
474 values of S_Y (ΔS_Y) and of E_{obs} (ΔE_{obs}) are depicted in Fig. 2e-3e and 2f3f,
475 respectively.

476 Values of E_Y , S_Y , and E_{obs} obtained at the end of the iteration procedure
477 through iES_ROM generally decrease with n . When $n = 25$ or 30 , the values of E_Y
478 and S_Y based on iES_ROM tend to approach their counterparts obtained through
479 iES_FSM. The latter generally correspond to the lowest values across TCs 1-6. These
480 findings are consistent with the observation (Xia et al., 2020; 2025) that accuracy of
481 ROM for $n = 30$ and FSM are very similar for solute transport. They are also in line

Formatted: Strikethrough

482 with the results of Li et al. (2013b), who documented a high degree of correlation
483 between simulated concentrations provided by their ROM and FSM for non-reactive
484 transport. ~~As expected, the solution accuracy of ROM increases with n . When n
485 approaches the dimension of the FSM (i.e., the representativeness of the basis
486 functions is essentially guaranteed), the ROM-based solution approaches its
487 FSM-based counterpart.~~

488 Figure 3-4 depicts E_Y (Fig. 3a4a), S_Y (Fig. 3b4b), and E_{obs} (Fig. 3e4c) for
489 TCs 7-12 obtained through iES_ROM and iES_FSM when the well characteristics are
490 deterministically known. Similar to above, results obtained through iES_FSM are
491 identical among TCs 7-12 and are taken as reference. Values of ΔE_Y , ΔS_Y , and
492 ΔE_{obs} are depicted in Fig. 3d4d, 3e4e, and 3f4f, respectively.

493 Consistent with what one can observe in Fig. 23, values of E_Y , S_Y , and E_{obs}
494 obtained at the end of the iteration procedure for TCs 7-12 through iES_ROM
495 generally decrease with n . Except for the cases where $n = 5$ or 10 (corresponding to
496 low solution accuracy of ROM), values of E_Y , S_Y , and E_{obs} for TCs 9-12 based on
497 either iES_ROM or iES_FSM are lower than their counterparts related to TCs 3-6.
498 These results suggest that the accuracy of conductivity estimates is lower when q_s is
499 uncertain compared to the case where q_s is deterministic.

500 Figure 4-5 depicts the values of E_{x_1} (Fig. 4a5a), E_{x_2} (Fig. 4b5b), E_{q_s} (Fig.
501 4e5c), S_{x_1} (Fig. 4d5d), S_{x_2} (Fig. 4e5e), and S_{q_s} (Fig. 4f5f) versus the number of
502 outer iterations for TCs 1-6 obtained through iES_ROM and iES_FSM. Values of E_{x_1} ,
503 E_{x_2} , and E_{q_s} obtained through iES_ROM approach their iES_FSM-based

504 counterparts as n increases. This is consistent with the observation that increasing n
505 improves the accuracy of the ROM-based solution (see also Li et al., 2013b),
506 therefore enhancing the accuracy of the identification of the well attributes.

507 Figure 5-6 depicts the estimated (ensemble) Y fields for TCs 1-6 obtained
508 through iES_ROM and iES_FSM, together with their reference Y field. The white
509 circle and cross symbols in Fig. 5-6 denote the estimated and reference locations of
510 the pumping well, respectively. As n increases, the estimated Y field obtained through
511 iES_ROM (Fig. 5a-5f) approaches its iES_FSM-based counterpart and the
512 reference Y field (Fig. 5h). The accuracy of the iES_ROM-based estimate of the
513 location of the pumping well generally increases with n , consistent with the nature of
514 the findings illustrated in Fig. 45. Figure 6-7 depicts the estimated (ensemble) Y
515 variance fields for TCs 1-6 based on iES_ROM and iES_FSM. The white circle and
516 cross symbols therein denote the identified and reference locations of the pumping
517 well, respectively. These results show that the variance of Y is overestimated when n
518 is small. This is related to the observation that small values of n correspond to large
519 modeling errors (i.e., low solution accuracy) of ROM (as also seen in Li et al. (2013b)
520 and Pasetto et al. (2017)). The latter, in turn, imprint the low accuracy of conductivity
521 estimates (see Fig. 5a-6a in the case of $n = 5$) and yield overestimated values for the
522 variance of Y (see Fig. 6a7a).

523 Figure 7-8 depicts the empirical probability density function (PDF) of x_{1,q_s} ,
524 x_{2,q_s} , and $\ln q_s$ at the end of the iteration procedure for TCs 1, 2, 4, and 6 as
525 obtained through iES_ROM and iES_FSM, together with their counterparts associated

526 with initial guess (black solid) and reference values (black dashed). One can observe
527 that large values of n yield high accuracy for x_{1,q_s} and x_{2,q_s} estimates, as visually
528 indicated by the compact supports associated with the empirical PDFs of x_{1,q_s} (Fig.
529 [7a8a](#)) and x_{2,q_s} (Fig. [7b8b](#)). The accuracy of the estimate of q_s is already
530 acceptable when $n = 5$.

531 As an additional element, we explore the way the choice of the value of n
532 impacts the local PDFs of hydraulic head and solute concentration. We do so upon
533 considering the results associated with three reference points (i.e., I, II, and III in Fig.
534 [4d2d](#)) that are aligned in the direction of the mean groundwater flow. Figure [8-9](#)
535 depicts the (sample) PDFs of (hydraulic) head at these three selected locations (Figs.
536 [8a9a-8e9c](#)) obtained through iES_ROM and iES_FSM at the end of the iteration
537 procedure for TCs 1, 2, 4, and 6. Black solid lines included therein indicate reference
538 head values. Note that the PDFs stemming from iES_FSM peak at values very close
539 to their reference counterparts. Hence, the corresponding empirical PDFs are
540 considered as reference. The logarithm absolute difference (Δ PDF, evaluated as the
541 pointwise log-ratio of the densities and corresponding to a local measure of relative
542 likelihood between two empirical PDFs) between the PDFs of the head at points I-III
543 obtained through iES_ROM based on diverse values of n and their counterpart based
544 on iES_FSM are also shown in Figs. [8d9d-8f9f](#), respectively. One can see that a large
545 value of n (e.g., $n = 30$ for TC6) corresponds to high accuracy of the PDF of head, as
546 quantified through a low value of Δ PDF. Although the head solution is obtained by
547 solving FSM, the accuracy of the conductivity estimate is impacted by n . The latter,

548 therefore, impacts the accuracy of heads. Fig. 9–10 depicts results related to solute
549 concentration. As expected, the PDFs stemming from iES_FSM peak at values very
550 close to their reference counterparts also in this case. Consistent with Fig. 89, a large
551 value of n (e.g., 30 for TC6) corresponds to high accuracy in the delineation of the
552 PDF of solute concentration.

553 As a complement to these results, values of the Kullback-Leibler Divergence
554 (KLD) between the (sample) PDFs of head at the three reference points at the last
555 outer iteration obtained through iES_FSM (h_{FSM}) and iES_ROM (h_{ROM}) with $n = 5$
556 (TC1), 10 (TC2), 20 (TC4), and 30 (TC6) are listed in Table S1 (see supplementary
557 information). We recall that values of $\text{KLD}(h_{\text{ROM}}||h_{\text{FSM}})$ (or $\text{KLD}(h_{\text{FSM}}||h_{\text{ROM}})$)
558 quantify (in a global sense) information loss when using h_{FSM} (h_{ROM}) to approximate
559 h_{ROM} (h_{FSM}). Values of $\text{KLD}(h_{\text{ROM}}||h_{\text{FSM}})$ generally increase with n . This indicates that
560 the difference between PDFs of h_{ROM} and h_{FSM} decrease as n increases. While the
561 highest values of $\text{KLD}(h_{\text{FSM}}||h_{\text{ROM}})$ correspond to $n = 5$, no clear decreasing trends
562 with increasing n are observed. Furthermore, the difference between $\text{KLD}(h_{\text{ROM}}||h_{\text{FSM}})$
563 and $\text{KLD}(h_{\text{FSM}}||h_{\text{ROM}})$ generally decreases as n increases. This is related to the
564 observation that the accuracy of ROM tends to increase as the dimension of the
565 reduced-order model increase. Values of KLD between the empirical PDFs of solute
566 concentrations at the three selected reference points at the last outer iteration obtained
567 through iES_FSM (c_{FSM}) and iES_ROM (c_{ROM}) with $n = 5$ (TC1), 10 (TC2), 20 (TC4),
568 and 30 (TC6) are listed in Table S2 (see supplementary information).

569 4.2 Effect of the ensemble size (Group B)

570 Figure 10-11 depicts iES_ROM- and iES_FSM-based values of E_Y (Fig.
571 10a11a), S_Y (Fig. 10b11b), and E_{obs} (Fig. 10e11c) versus the number of outer
572 iterations for TCs 6 and 13-15. Values of E_Y and E_{obs} decrease as the ensemble
573 size N_{MC} increases (while the value of S_Y increases) regardless of the approach
574 employed. With reference to TC13, we note that when $N_{MC} = 30$ the values of E_Y
575 decrease during the course of the first outer iterations to then increase during the last
576 outer iterations, values of S_Y dropping rapidly during the iteration procedure,
577 regardless of the approach employed. This phenomenon is typically linked to the
578 occurrence of filter inbreeding caused by a limited ensemble size (Chen and Zhang,
579 2006; Xia et al., 2018; 2024). Values of E_Y and S_Y for TCs 6 and 13-15 obtained
580 through iES_ROM are overall similar to those associated with iES_FSM. The
581 iES_ROM-based value of E_{obs} obtained at the end of the iteration procedure for a
582 given TC is typically larger than its iES_FSM-based counterpart. This is linked to the
583 observation that the limited system dimension of ROM induces low accuracy of
584 concentrations and (possibly) heads due to low accuracy of conductivity estimates,
585 pumping rate, and well locations.

586 Figure 11-12 depicts the values of E_{x_1} (Fig. 11a12a), E_{x_2} (Fig. 11b12b), E_{q_s}
587 (Fig. 11e12c), S_{x_1} (Fig. 11d12d), S_{x_2} (Fig. 11e12e), and S_{q_s} (Fig. 11f12f) versus
588 the number of outer iterations for TCs 6 and 13-15 obtained through iES_ROM and
589 iES_FSM. When increasing N_{MC} , values of E_{x_1} , E_{x_2} , and E_{q_s} obtained through
590 either iES_ROM or iES_FSM do not show a clear trend. Values of S_{x_1} , S_{x_2} , and S_{q_s}

591 generally increase with N_{MC} , a result that is consistent with the findings encapsulated
592 in Fig. [40b11b](#). Similar findings are also documented by Xu and Gómez-Hernández
593 (2018, their Fig. 17), who show that, when considering joint identification of
594 contaminant sources and hydraulic conductivities, the accuracy of estimates of key
595 attributes characterizing contaminant sources does not necessarily improve after some
596 time and as data assimilation progresses. We further note that jointly estimating
597 conductivity and identifying source/sink term attributes (in terms of flow rate and
598 location) is associated with a highly nonlinear optimization process. Hence, the
599 accuracies of location and pumping rate estimation through iES_FSM are not always
600 higher than those stemming from iES_ROM in terms of the values of the metrics
601 employed (i.e., E_{x_1} , E_{x_2} , and E_{q_s}).

602 Figures [42-13](#) depicts the estimated (ensemble mean) Y fields for TCs 6 and
603 13-15 obtained through iES_ROM and iES_FSM. Figure [43-14](#) depicts the associated
604 Y variance fields for TCs 6 and 13-15 obtained through iES_ROM and iES_FSM. The
605 white (black) circle and cross symbols in Fig. [42-13](#) (or Fig. [43-14](#)) represent the
606 identified and the reference locations of the pumping well, respectively. Visual
607 comparison of Fig. [42-13](#) and Fig. [5h-6h](#) suggests that the estimated Y fields rendered
608 through an ensemble size $N_{MC} = 100$ (i.e., TC14) obtained through iES_ROM and
609 iES_FSM are the closest ones to the reference Y field. Nevertheless, jointly analyzing
610 Figs. [40a11a](#), [4213](#), and [43-14](#) reveal that the estimated Y field corresponding to N_{MC}
611 = 10,000 (TC6) obtained through iES_ROM is the one most closest to the reference Y
612 field in terms of E_Y ($= 0.41$). Additionally, the identified and reference locations of the

613 pumping well obtained through either iES_ROM or iES_FSM are close to each other,
614 thus supporting the capability of both approaches to identify the well location.

615 **4.3 Effect of quality and available number of observations (Group C)**

616 Table 2 lists values of E_Y , S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s} at
617 the end of iteration procedure for TC16 (characterized by $\sigma_{obs} = 0.001$), TC6 (σ_{obs}
618 $= 0.01$), and TC17 ($\sigma_{obs} = 0.1$) obtained through iES_ROM and iES_FSM. Values of
619 E_Y , S_Y , E_{obs} , E_{x_1} , and E_{q_s} generally increase as the quality of observations
620 deteriorates, i.e., σ_{obs} increasing from 0.001 to 0.1. These results are also consistent
621 with prior findings by Xia et al. (2018) according to which accuracy of conductivity
622 estimates increases as the quality of observations improves. Values of E_{x_2} obtained
623 through iES_ROM and iES_FSM do not monotonically decrease as σ_{obs} decreases.
624 This is typically related to the strong nonlinear nature associated with the optimization
625 process (see also Xu and Gómez-Hernández, 2018).

626 Figure 14-15 depicts iES_ROM- and iES_FSM-based values of E_Y (Fig.
627 14a15a), S_Y (Fig. 14b15b), and E_{obs} (Fig. 14e15c) versus the number of outer
628 iterations for TCs 6 and 18-19. Values of E_Y (or S_Y) for TCs 18 (where the number
629 of monitoring wells is $N_m = 9$), 19 ($N_m = 18$), and 6 ($N_m = 55$) obtained through
630 iES_ROM are similar to their iES_FSM-based counterparts and decrease as N_m
631 increases. Values of E_{obs} obtained through iES_FSM decrease as N_m increases,
632 while iES_ROM-based results do not display a clear trend with N_m . This result may
633 be attributed to the fact that, while increasing the number of monitoring wells
634 enhances the amount of information available for estimating hydraulic conductivity,

635 errors introduced through model reduction influence the evolution of the solute
636 concentration mismatch between observations and simulations during the iterative
637 calibration process.

638 Figure 15-16 depicts the values of E_{x_1} (Fig. 15a16a), E_{x_2} (Fig. 15b16b), E_{q_s}
639 (Fig. 15e16c), S_{x_1} (Fig. 15d16d), S_{x_2} (Fig. 15e16e), and S_{q_s} (Fig. 15f16f) versus
640 the number of outer iterations for TCs 6 and 18-19 obtained through iES_ROM and
641 iES_FSM. Values of E_{x_1} (E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , S_{q_s} , or S_Y) for TCs 18 (for a
642 number $N_m = 9$ of monitoring wells), 19 ($N_m = 18$), and 6 ($N_m = 55$) obtained
643 through either iES_ROM or iES_FSM decrease as N_m increases. Values of the same
644 metric (i.e., E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , or S_{q_s}) obtained through iES_ROM and
645 iES_FSM are overall close to each other.

646 Figure 16-17 depicts the empirical PDF of x_{1,q_s} , x_{2,q_s} , and $\ln q_s$ at the end of
647 the iteration procedure for TCs 18, 19, and 6 obtained through iES_ROM and
648 iES_FSM, together with their reference counterparts (black dashed lines). One can
649 observe that increasing N_m leads to improved accuracy of the identification of
650 pumping well attributes, as suggested by the reduced support and location of the
651 peaks of the PDFs of x_{1,q_s} (Fig. 16a17a), x_{2,q_s} (Fig. 16b17b), and $\ln q_s$ (Fig.
652 16c17c) obtained through either iES_ROM or iES_FSM and observed as N_m varies
653 from 9 to 55. On the basis of these results, it is hard to tell which approach provides
654 higher accuracy of pumping well identification, solely in terms of Fig. 1617. To
655 complement these findings, Table S3 (see supplementary information) lists the values
656 of KLD between the empirical PDFs of x_{1,q_s} (x_{2,q_s} , or $\ln q_s$) obtained through

657 iES_FSM (denoted as p_{FSM}) and iES_ROM (denoted as p_{ROM}) with $n = 30$,
 658 considering $N_m = 9$ (TC18), 18 (TC19), and 55 (TC6), respectively. Values of
 659 $\text{KLD}(p_{\text{ROM}}||p_{\text{FSM}})$ (with $j = x_{1,q_s}, x_{2,q_s}$, and $\ln q_s$ show an overall decreasing trend
 660 as N_m increase, while $\text{KLD}(p_{\text{FSM}}||p_{\text{ROM}})$ consistently decreases with N_m . These
 661 results are consistent with the observation that increasing the number of monitoring
 662 wells improves the accuracy of conductivity estimates (as also seen by Tong et al.
 663 (2010) and Xia et al. (2018)) as well as pumping rate and well location through both
 664 approaches, thus, in turn, reducing discrepancies between the corresponding PDFs.

665 **4.4 Effect of the mean and variance of the initial ensemble of Y (Group D)**

666 Table 3 lists the values of $E_Y, S_Y, E_{obs}, E_{x_1}, E_{x_2}, E_{q_s}, S_{x_1}, S_{x_2}$, and S_{q_s}
 667 at the end of the iteration procedure for TCs 20 (characterized by a mean $\mu = -0.5$ of
 668 the initial ensemble of Y), 6 ($\mu = 1.2$), and 21 ($\mu = 2.0$) obtained through
 669 iES_ROM and iES_FSM. We recall that the mean value employed to generate the
 670 reference Y field is equal to 0.8. When the discrepancy between μ and the mean
 671 value of the reference Y field increases, the error metrics employed display an overall
 672 increase, E_{x_1} and E_{q_s} constituting notable exceptions. This finding is consistent
 673 with the behavior documented by Xia et al. (2024) who considered two
 674 correlation-based localization approaches to assess conductivity estimation accuracy
 675 with respect to the mean of the initial ensemble of Y . ~~Overall, this result is related to~~
 676 ~~the highly nonlinear nature of the optimization process (see also Sections 4.2 and 4.3)~~
 677 ~~stemming from the joint estimation of conductivities and pumping well attributes.~~

678 Table 4 lists the values of $E_Y, S_Y, E_{obs}, E_{x_1}, E_{x_2}, E_{q_s}, S_{x_1}, S_{x_2}$, and S_{q_s}

Formatted: Strikethrough

679 at the end of the iteration procedure for TCs 22 (characterized by a variance $\sigma_Y^2 =$
680 0.01 of the initial ensemble of Y), 6 ($\sigma_Y^2 = 1.0$), and 23 ($\sigma_Y^2 = 2.0$) obtained through
681 iES_ROM and iES_FSM. We recall that the reference Y field is characterized by a
682 unit variance. The values of E_Y and E_{x_2} obtained through both approaches increase
683 as the discrepancy between σ_Y^2 and the variance of the reference Y field increases.
684 The values of E_{obs} , E_{x_1} , and E_{q_s} obtained through both approaches generally
685 increase with σ_Y^2 . Similarly, values of metrics employed to quantify variability of the
686 final ensemble of realizations (i.e., S_Y , S_{x_1} , S_{x_2} , and S_{q_s}) consistently increase
687 with σ_Y^2 .

688 A joint analysis of the results illustrated in Sections 4.1, 4.2, and 4.3 suggests
689 that E_Y and S_Y provided by both approaches show a consistent behavior as a
690 function of the key feature of interest. Otherwise, the response of the metrics
691 associated with the pumping well attributes provided by both approaches reflects the
692 enhanced nonlinearity of the associated optimization process. Additionally, the
693 accuracy of the conductivity estimate possibly contributes more to the minimization
694 of the objective function than that of pumping well identification. Additionally, the
695 values of the metrics in Sections 4.1, 4.2, and 4.3 provided by the two approaches are
696 generally consistent with each other, thus supporting the representativeness of the
697 iES_ROM-based results.

698

699 **4.5 Effect of the snapshot size (Group E)**

700 Table 5 lists percentage differences of the values of the performance metrics

701 considered (i.e., E_Y , S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s}) at the end of
702 the iteration procedure for TCs 24-28 obtained through iES_ROM, considering their
703 counterparts through TC6 as references. These results show that the values of E_Y
704 and S_Y systematically decrease as N_{sn} increases from 30 to 1,000, while the other
705 metrics display an overall decreasing pattern. This is related to the observation that a
706 larger snapshot size corresponds to a higher accuracy of basis functions (Pasetto et al.,
707 2014). Otherwise, it is worth noting that snapshots are evaluated only once throughout
708 the entire data assimilation processes, thus resulting in a limited computational cost.

709 A CPU time of about 13 minutes is required for running TC28 (using a processor
710 13th Gen Intel(R) Core(TM) i7-13700K 3.40 GHz with 32 GB RAM). The CPU time
711 required to complete TC6 upon relying on iES_FSM (122 minutes) is about 9 times
712 the corresponding CPU time required to complete TC28 through iES_ROM (28
713 minutes), percentage differences associated with E_Y and S_Y being equal to 0.50%
714 and 0.21%, respectively. The CPU required time for running TC28 is about 13 minutes
715 (using a processor 13th Gen Intel(R) Core(TM) i7-13700K 3.40 GHz with 32 GB
716 RAM). CPU times for running TC6 through iES_ROM and iES_FSM are about 28
717 and 122 minutes, respectively. The CPU time to complete TC6 upon relying on
718 iES_FSM is about 9 times the CPU time for TC28 through iES_ROM, while
719 percentage differences associated with E_Y and S_Y are 0.50% and 0.21%,
720 respectively. CPU time savings can become more pronounced during data
721 assimilation for a groundwater system of large size, due to the higher memory
722 requirements of iES_FSM for storing and computing large-dimensional vectors and

Formatted: Font color: Auto

Field Code Changed

Formatted: Font color: Auto

Field Code Changed

Formatted: Font color: Auto

723 matrices as compared to iES_ROM. ~~These results support the ability of iES_ROM to~~
724 ~~estimate conductivity and identify the uncertain features of a pumping well under the~~
725 ~~conditions analyzed.~~

Formatted: Strikethrough

726 Additionally, we emphasize that relying on realizations of Y associated with
727 (spatial) statistics different from their theoretical counterparts linked to the initial
728 ensemble of Y fields can contribute to deteriorate the quality of the selected snapshots.
729 Low quality snapshots yield low quality basis functions and low accuracy of ROM
730 outcomes (see our results in Section 4.1; Pasetto et al., 2014; Xia et al., 2020). The
731 latter deteriorate the accuracy of conductivity estimates and pumping well attributes.
732 Additional studies should be devoted to assess the potential of techniques (including,
733 e.g., greedy algorithms) that might contribute to increase the quality of snapshots.

Formatted: Font color: Auto

734 5. Conclusions

735 This study addresses ~~ed the~~ joint estimation of (uncertain, spatially heterogeneous)
736 hydraulic conductivities and attributes (location and flow rate) of a pumping well in a
737 two-dimensional confined aquifer in the presence of (non-reactive) solute transport
738 taking place across a steady-state flow field. Our analyses rest on an iterative
739 Ensemble Smoother (iES) coupled with a Reduced-Order Model (ROM) for solute
740 transport (the overall strategy being denoted as iES_ROM). The ROM is constructed
741 via Proper Orthogonal Decomposition (POD), using basis functions derived from the
742 numerical solutions of the Full System Model (FSM) over the entire simulation period.
743 The pumping well is characterized by its spatial coordinates (x_{1,q_s}, x_{2,q_s}) and a
744 constant pumping rate q_s . The ROM can achieve a solution accuracy similar to that

745 of the FSM, while significantly reducing computational demands. Notably, as stated
746 above, the basis functions are computed only once throughout the iES_ROM iteration
747 process, thus further enhancing efficiency. As a benchmark, the traditional iES
748 approach relying on the FSM (termed iES_FSM) is also implemented to estimate
749 conductivity and identify well attributes.

750 To assess the performance and robustness of the proposed iES_ROM approach,
751 twenty-eight test cases (TCs 1-28) are designed and structured according to five
752 categories (Groups A-E; Section 3), each targeting different influencing factors. These
753 include the dimension of the reduced-order model (n), ensemble size (N_{mc}), standard
754 deviation of the white noise representing measurement error (σ_{obs}), number of
755 monitoring wells (N_m), mean (μ) and variance (σ_y^2) of the initial log-conductivity
756 field, and snapshot size (N_{sn}). The performance of iES_ROM is systematically
757 compared with that of iES_FSM using nine evaluation metrics, encompassing the
758 absolute error (E_Y ; Equation (9)) and estimated standard deviation (S_Y ; Equation (9))
759 between estimated and reference values of Y ; the absolute errors and estimated
760 standard deviations of the pumping well coordinates and rate (Equation (10)); and the
761 average absolute difference between simulated and reference observations (E_{obs} ;
762 Equation (12)).

763 Our work leads to the following major conclusions.

- 764 1. Both iES_ROM and iES_FSM yield accurate estimates of hydraulic
765 conductivity distributions and identify the pumping well attributes across a
766 wide range of tested conditions, including variations in model dimension,

767 ensemble size, measurement noise, number of monitoring wells, and
768 statistical properties of the initial ensemble.

769 2. The iES_ROM approach achieves estimation accuracy similar to that of
770 iES_FSM when using a moderate reduced-order dimension ($n = 25$ or 30).
771 Otherwise, relying on a small dimension (e.g., $n = 5$) yields filter divergence
772 due to unaccounted model errors. Increasing n effectively mitigates this
773 issue and enhances the stability of the iES_ROM performance.

774 3. When hydraulic conductivity and pumping well attributes are jointly
775 estimated, both iES_ROM and iES_FSM exhibit a slight reduction in the
776 accuracy of conductivity estimates compared to scenarios where only
777 conductivity is estimated. This trend is reflected in the values of E_Y , S_Y ,
778 and E_{obs} across TCs 1-12. Under such joint estimation, results in terms of
779 E_Y , S_Y , and E_{obs} with respect to different influencing factors remain of
780 acceptable quality for both iES_ROM and iES_FSM, consistent with the
781 patterns observed in conductivity-only estimation. The behaviors of the
782 remaining performance metrics are mutually consistent and within
783 acceptable ranges, although somewhat less orderly.

784 4. Relying on the iES_ROM approach yields an accuracy similar to that of
785 iES_FSM in estimating hydraulic conductivity and identifying pumping well
786 attributes for both moderate ($N_{sn} = 500$ or 1000) and large ($N_{sn} = 10000$)
787 ensemble sizes. This result supports its robustness with respect to ensemble
788 size selection.

789 5. In terms of computational efficiency, iES_ROM yields substantial time
790 savings compared to iES_FSM. For instance, with $N_{sm} = 500$ and $n = 30$,
791 the CPU times for iES_ROM and iES_FSM are approximately 28 and 122
792 minutes, respectively (i.e., iES_FSM requires a computation time that is
793 about nine times longer while yielding similar estimation accuracy).

794 Additional elements of interest associated with future studies on coupling iES
795 with ROM include the analysis of transient saturated/unsaturated flow, reactive
796 transport, and density-dependent flow/transport scenarios. When considering
797 nonlinear systems, reliance on discrete matrix interpolation schemes (Negri et al.,
798 2015; Bonomi et al., 2017) constitutes a promising approach to enhance
799 computational advantages of ROM.

800 5. Moreover, the values of N_{MC} that one should consider in a field application are
801 case-dependent. In this context, localization techniques can be embedded in DA
802 processes, as these can reduce negative influences of spurious correlation on
803 parameter estimate arising from reliance on small ensemble sizes.

804 *Author contributions.* All authors contributed to the preparation of the manuscript.

805 *Acknowledgments.* This work was supported by the National Nature Science
806 Foundation of China (Grant No. 42002247), Nature Science Foundation of Fujian
807 Province, China (Grant No. 2025J01529; 2025J08248), and Opening Fund of Key
808 Laboratory of Geohazard Prevention of Hilly Mountains, Ministry of Natural
809 Resources (FJKLGH2024K008). M.R. acknowledges funding from the National
810 Recovery and Resilience Plan (NRRP), mission 4 component 2 investment 1.4 - call

Formatted: Font color: Auto

Formatted: Indent: First line: 2 ch, No bullets or numbering

Formatted: Font color: Auto

Formatted: Font color: Auto

811 for tender no. 3138 of 16 December 2021, rectified by decree no. 3175 of 18
812 December 2021 of Italian Ministry of University and Research funded by the
813 European Union - NextGenerationEU, project code CN_00000033, concession decree
814 no. 1034 of 17 June 2022 adopted by the Italian Ministry of University and Research,
815 CUP D43C22001250001, project title “National Biodiversity Future Center - NBFC”.

816 **References**

817 Asher, M.J., Croke, B.F.W., Jakeman, A.J., Peeters, L.J.M., 2015. A review of
818 surrogate models and their application to groundwater modeling. *Water Resour.*
819 *Res.* 51, 5957–5973.

820 Ballio, F., Guadagnini, A., 2004. Convergence assessment of numerical Monte Carlo
821 simulations in groundwater hydrology. *Water Resour. Res.* 40, W04603.

822 Boyce, S.E., Nishikawa, T., Yeh, W.W.G., 2015. Reduced order modeling of the
823 newton formulation of modflow to solve unconfined groundwater flow. *Adv.*
824 *Water Resour.* 83, 250–262.

825 [Bonomi, D., Manzoni, A., Quarteroni, A., 2017. A matrix DEIM technique for model](#)
826 [reduction of nonlinear parametrized problems in cardiac mechanics. *Comput.*](#)
827 [Methods Appl. Mech. Eng.](#) 324, 300-326.

828 Chen, Y., Zhang, D., 2006. Data assimilation for transient flow in geologic formations
829 via ensemble Kalman filter. *Adv Water Resour.*, 29(8): 1107-22.

830 Chen, Y., Oliver, D. S., 2013. Levenberg–Marquardt forms of the iterative ensemble
831 smoother for efficient history matching and uncertainty quantification.
832 *Computational Geosciences*, 17(4): 689-703.

833 Chen, Z., Jaime Gomez-Hernandez, J., Xu, T., Zanini, A., 2018. Joint identification of
834 contaminant source and aquifer geometry in a sandbox experiment with the
835 restart ensemble kalman filter. *J. Hydrol.* 564, 1074–1084.

836 Deutsch, C.V., Journel, A.G., 1998. *GSLIB: Geostatistical Software Library and*
837 *User's Guide*, second ed. Oxford University Press, New York.

838 Evensen, G., 2009. *Data Assimilation: The Ensemble Kalman Filter. Data*
839 *Assimilation: The Ensemble Kalman Filter.*

840 Ju, L., Zhang, J., Meng, L., et al. 2018. An adaptive Gaussian process-based iterative
841 ensemble smoother for data assimilation. *Adv. Water Resour.*, 115: 125-35.

842 Li, X., Chen, X., Hu, B.X., Navon, I.M., 2013a. Model reduction of a coupled
843 numerical model using proper orthogonal decomposition. *J. Hydrol.* 507,
844 227–240.

845 Li, X., Hu, B. X., 2013b. Proper orthogonal decomposition reduced model for mass
846 transport in heterogeneous media. *Stochastic Environmental Research and Risk*
847 *Assessment*, 27(5): 1181-91.

848 Luo, Z., Li, H., Zhou, Y., et al. 2012. A reduced finite element formulation based on
849 POD method for two-dimensional solute transport problems. *Journal of*
850 *Mathematical Analysis and Applications*, 385(1): 371-83.

851 Luo, X., Bhakta, T., 2020. Automatic and adaptive localization for ensemble-based
852 history matching. *Journal of Petroleum Science and Engineering*, 184: 106559.

853 Mo, S., Zabarar, N., Shi, X., et al. 2019. Deep Autoregressive Neural Networks for
854 High-Dimensional Inverse Problems in Groundwater Contaminant Source

855 Identification. Water Resour. Res., 55(5): 3856-81.

856 [Negri, F., Manzoni, A., Amsallem, D., 2015. Efficient model reduction of](#)
857 [parametrized systems by matrix discrete empirical interpolation. J. Comput. Phys.](#)
858 [303, 431-454.](#)

859 Pasetto, D., Guadagnini, A., Putti, M., 2011. POD-based Monte Carlo approach for
860 the solution of regional scale groundwater flow driven by randomly distributed
861 recharge. Adv. Water Resour., 34(11): 1450-1463.
862 DOI:10.1016/j.advwatres.2011.07.003

863 Pasetto, D., Putti, M., Yeh, W.W.G., 2013. A reduced-order model for groundwater
864 flow equation with random hydraulic conductivity: Application to Monte Carlo
865 methods. Water Resour. Res., 49(6): 3215-3228. DOI:10.1002/wrcr.20136

866 Pasetto, D., Guadagnini, A., Putti, M., 2014. A reduced-order model for Monte Carlo
867 simulations of stochastic groundwater flow. Computational Geosciences, 18(2):
868 157-169. DOI:10.1007/s10596-013-9389-4

869 Pasetto, D., Ferronato, M., Putti, M., 2017. A reduced order model-based
870 preconditioner for the efficient solution of transient diffusion equations.
871 International Journal for Numerical Methods in Engineering, 109(8): 1159-1179.
872 DOI:10.1002/nme.5320

873 Pinnau, R., 2008. Model Reduction via Proper Orthogonal Decomposition /Schilders,
874 W. H. A., Van Der Vorst, H. A., Rommes, J. Model Order Reduction: Theory,
875 Research Aspects and Applications. Berlin, Heidelberg; Springer Berlin
876 Heidelberg, 95-109.

877 Razavi, S., Tolson, B.A., Burn, D.H., 2012. Review of surrogate modeling in water
878 resources. *Water Resour. Res.* 48

879 Rizzo, C., de Barros, F., Perotto, S., Oldani, L., Guadagnini, A., 2018. Adaptive POD
880 model reduction for solute transport in heterogeneous porous media. *Computat.*
881 *Geosci.* 22, 297–308.

882 Stanko, Z.P., Boyce, S.E., Yeh, W.W.-G., 2016. Nonlinear model reduction of
883 unconfined groundwater flow using pod and deim. *Adv. Water Resour.* 97,
884 130–143.

885 Tong, J., Hu, B. X., Yang, J., 2010. Using data assimilation method to calibrate a
886 heterogeneous conductivity field conditioning on transient flow test data.
887 *Stochastic Environmental Research and Risk Assessment*, 24(8): 1211-23.

888 Xia, C.-A., Luo, X., Hu, B.X., Riva, M., Guadagnini, A., 2021. Data assimilation with
889 multiple types of observation boreholes via the ensemble Kalman filter
890 embedded within stochastic moment equations. *Hydrol. Earth Syst. Sci.*, 25(4):
891 1689-1709. DOI:10.5194/hess-25-1689-2021

892 Xia, C.-A., Pasetto, D., Hu, B.X., Putti, M., Guadagnini, A., 2020. Integration of
893 moment equations in a reduced-order modeling strategy for Monte Carlo
894 simulations of groundwater flow. *J. Hydrol.*, 590: 125257.
895 DOI:<https://doi.org/10.1016/j.jhydrol.2020.125257>

896 Xia, C.-A., Guadagnini, A., Hu, B. X., Riva, M., Ackerer, P., 2019. Grid convergence
897 for numerical solutions of stochastic moment equations of groundwater flow,
898 *Stoch. Environ. Res. Risk Assess.*, 33 (8-9), 1565-1579,

899 <https://doi.org/10.1007/s00477-019-01719-6>.

900 Xia, C.-A., Hu, B.X., Tong, J., Guadagnini, A., 2018. Data Assimilation in
901 Density-Dependent Subsurface Flows via Localized Iterative Ensemble Kalman
902 Filter. *Water Resour. Res.*, 54(9): 6259-6281. DOI:10.1029/2017wr022369

903 Xia, C.-A., Li, J., Riva, M., et al., 2024. Characterization of conductivity fields
904 through iterative ensemble smoother and improved correlation-based adaptive
905 localization. *J. Hydrol.*, 634: 131054.

906 Xia, C.-A., Wang, H., Jian, W., et al., 2025. Reduced-order Monte Carlo simulation
907 framework for groundwater flow in randomly heterogeneous composite
908 transmissivity fields. *J. Hydrol.*, 651: 132593.

909 Xu, T., Jaime Gomez-Hernandez, J., 2018. Simultaneous identification of a
910 contaminant source and hydraulic conductivity via the restart normal-score
911 ensemble Kalman filter. *Adv. Water Resour.*, 112: 106-123.
912 DOI:10.1016/j.advwatres.2017.12.011

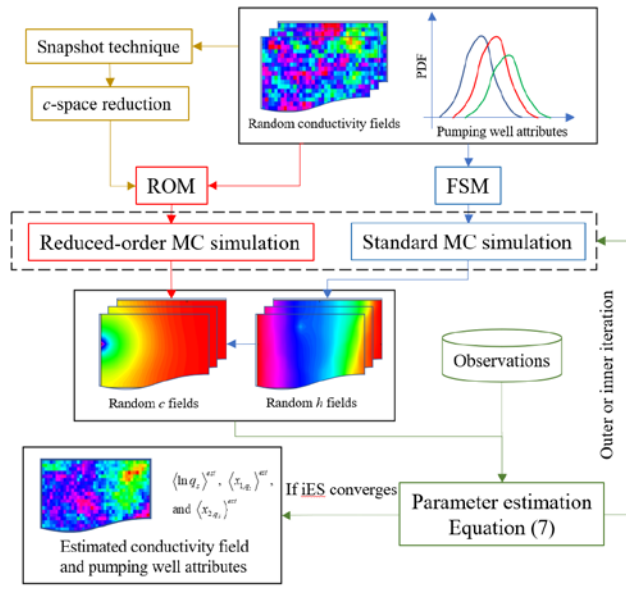
913 Zhang, D., 2002. *Stochastic Method for Flow in Porous Media – Coping with*
914 *Uncertainties*. Academic Press, Sand Diego, California.

915 Zhang J, Lin G, Li W, et al. An Iterative Local Updating Ensemble Smoother for
916 Estimation and Uncertainty Assessment of Hydrologic Model Parameters With
917 Multimodal Distributions. *Water Resour Res*, 2018, 54(3): 1716-33.

918

919
920
921
922
923
924

Figures



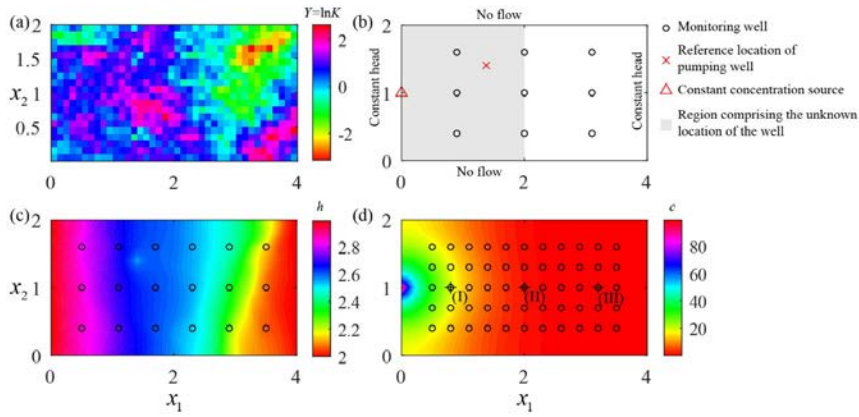
925

926 Fig. 1 Workflow of iES ROM, comprising (i) standard MC simulation of
 927 groundwater flow (relying on FSM), (ii) reduced-order MC approach for solute
 928 transport (relying on ROM), and (iii) iES coupled with ROM.

Formatted: Font: Italic
 Formatted: Font: Italic
 Formatted: Font: Italic

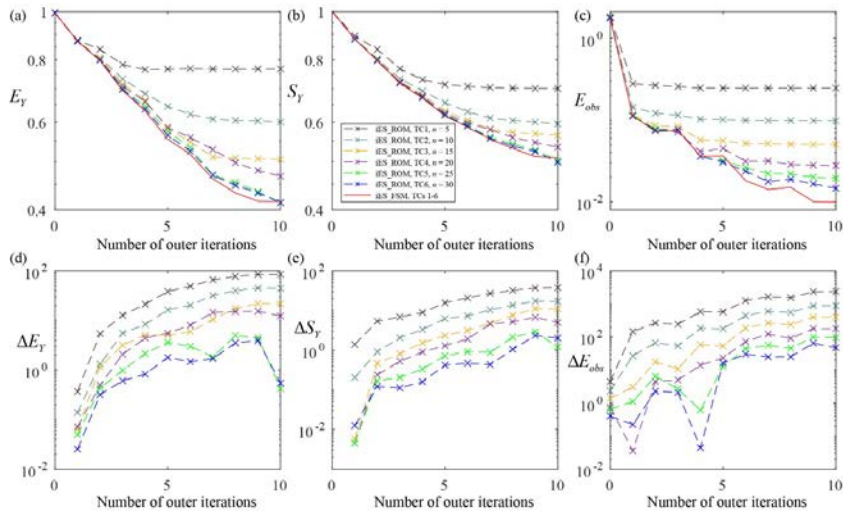
929
930

Formatted: Justified



931
932
933
934
935
936
937
938
939
940
941
942

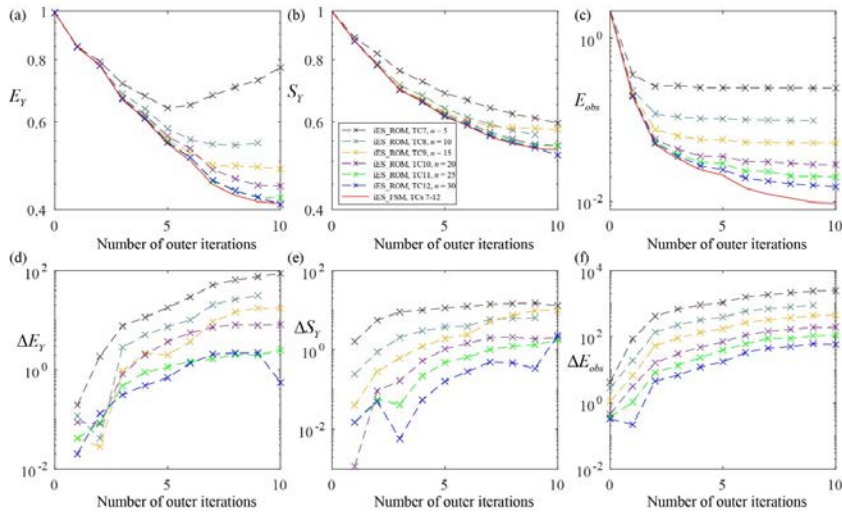
Fig. 4-2 (a) Reference field of $Y = \ln K$; (b) boundary conditions for groundwater flow and solute transport together with spatial distribution of 9 monitoring wells and reference location for the pumping well (shaded gray area corresponds to the region comprising the unknown location of the well); (c) hydraulic head corresponding to the reference Y field; and (d) solute concentration corresponding the reference Y field at final time step, including three selected locations (i.e., I, II, and III) at which empirical probability density functions of solute concentration is computed and considered for illustration purposes. Circles in (b), (c) and (d) correspond to the location of the 9, 18 and 55 monitoring wells, respectively, employed in the study (see Section 3 and Table 1).



944

945 Fig. 2-3 Values of (a) E_Y , (b) S_Y , and (c) E_{obs} versus the number of outer iterations
 946 obtained through iES_ROM considering various dimensions of reduced-order model
 947 (with $n = 5, 10, 15, 20, 25,$ and 30 for TCs 1-6, respectively) and iES_FSM (which
 948 provides identical results for TCs 1-6) for ensemble size $N_{MC} = 10,000$;
 949 corresponding percentage differences between the values of (d) E_Y (ΔE_Y), (e) S_Y
 950 (ΔS_Y), and (f) E_{obs} (ΔE_{obs}) evaluated through iES_ROM and iES_FSM.

951

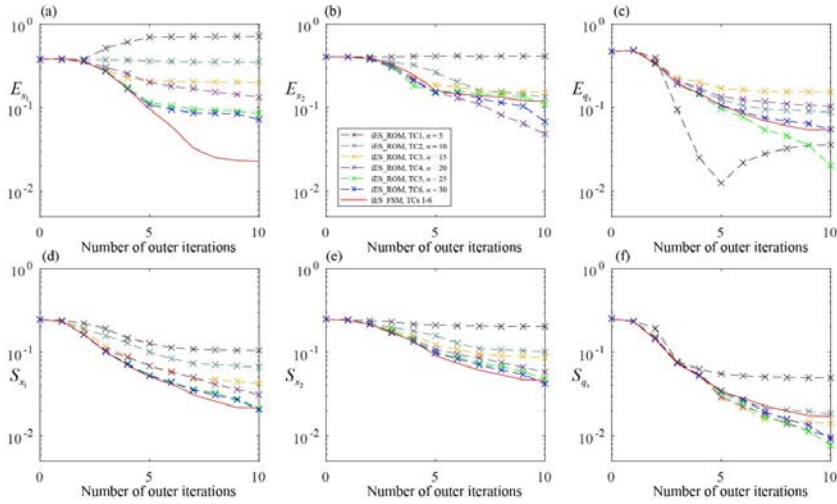


952

953 Fig. 3-4 Values of (a) E_Y , (b) S_Y , and (c) E_{obs} versus the number of outer iterations
954 obtained through iES_ROM considering various dimensions of reduced-order model
955 (with $n = 5, 10, 15, 20, 25,$ and 30 for TCs 7-12, respectively) and iES_FSM (which
956 provides identical results for TCs 7-12, when the pumping rate and location are
957 previously known and for an ensemble size $N_{MC} = 10,000$; corresponding percentage
958 differences between the values of (d) E_Y (ΔE_Y), (e) S_Y (ΔS_Y), and (f) E_{obs}
959 (ΔE_{obs}) evaluated through iES_ROM and iES_FSM.

960

961

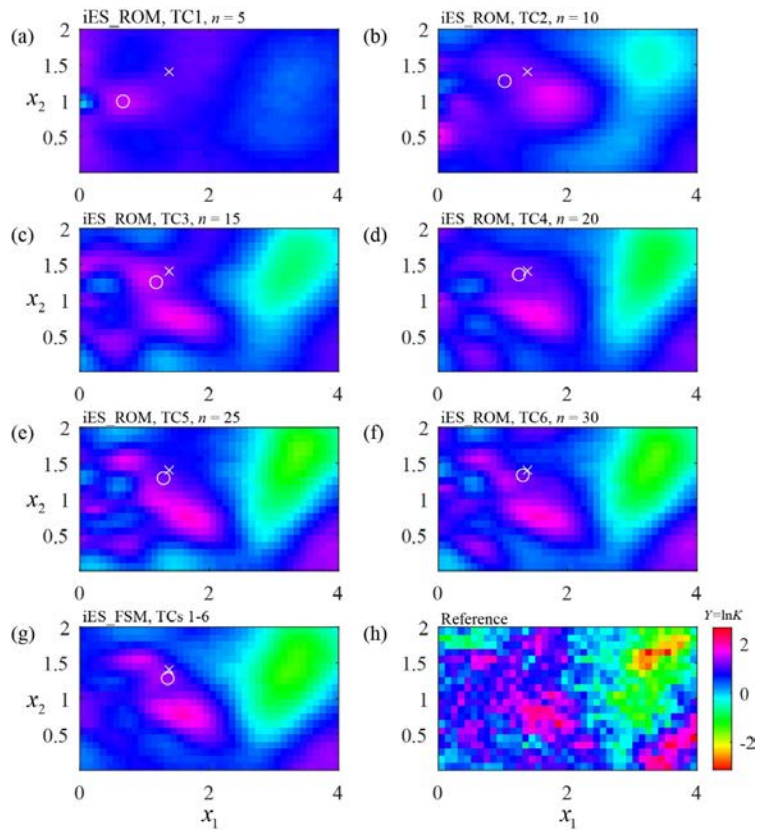


962

963 Fig. 4-5 Values of (a) E_{x_1} , (b) E_{x_2} , (c) E_{q_1} , (d) S_{x_1} , (e) S_{x_2} , and (f) S_{q_1} versus the

964 number of outer iterations obtained through iES_ROM considering various
 965 dimensions of reduced-order model (with $n = 5, 10, 15, 20, 25,$ and 30 for TCs 1-6,
 966 respectively) and iES_FSM (which provides identical results for TCs 1-6), when N_{MC}
 967 $= 10,000$.

968



969

970

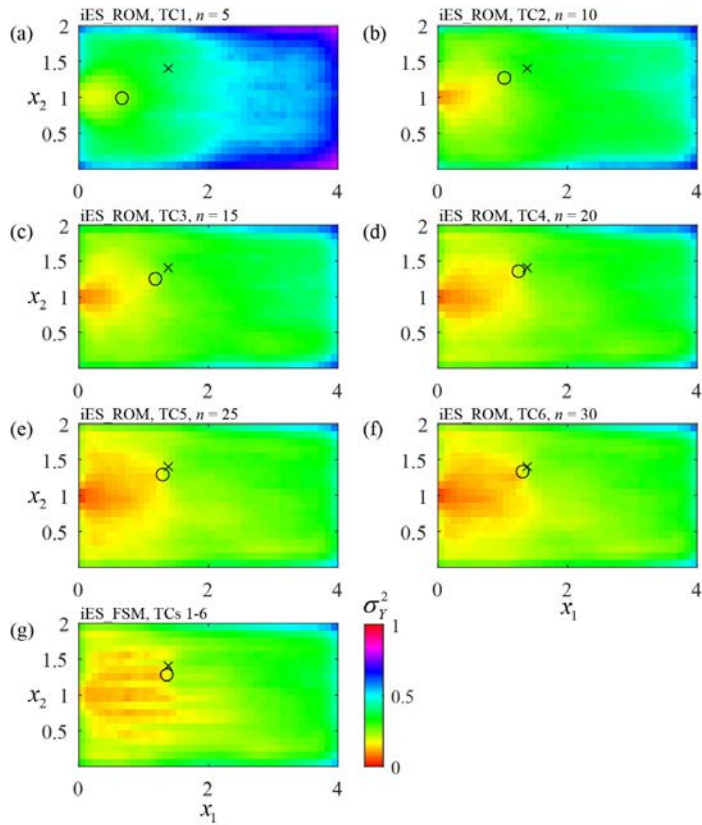
971

972

973

974

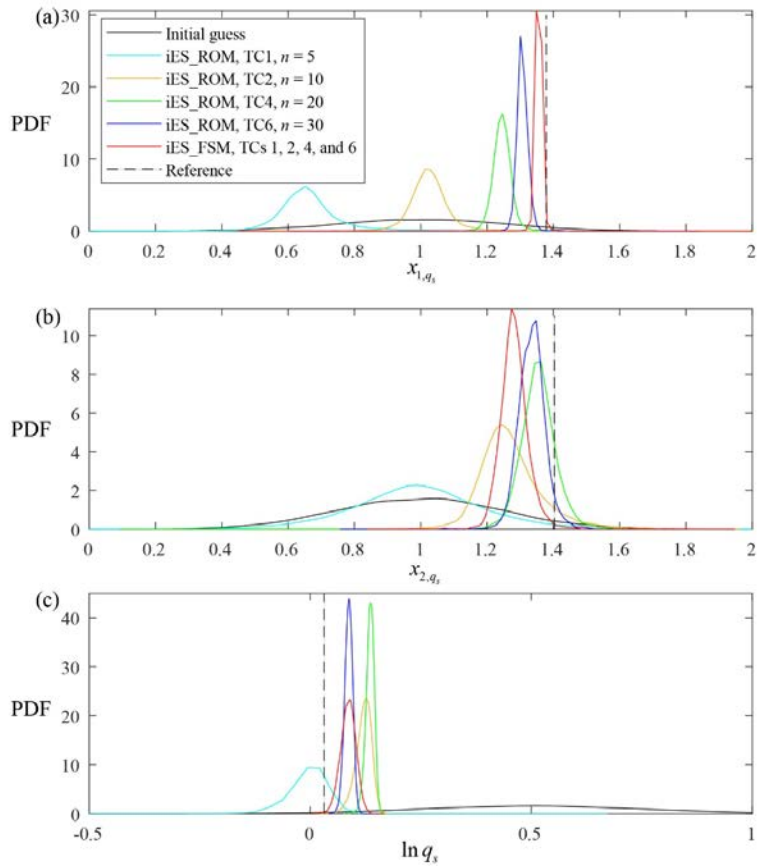
Fig. 5-6 Estimated (ensemble) mean Y fields at the final outer iteration through iES_ROM considering different n (equal to (a) 5, (b) 10, (c) 15, (d) 20, (e) 25, and (f) 30 for TCs 1-6, respectively) and (g) iES_FSM (which provides identical results for TCs 1-6) when $N_{MC} = 10,000$; (h) reference Y field.



976

977 Fig. 6-7 Estimated (ensemble) Y variance fields at the final outer iteration through
 978 iES_ROM considering different n (equal to (a) 5, (b) 10, (c) 15, (d) 20, (e) 25, and (f)
 979 30 for TCs 1-6, respectively) and (g) iES_FSM (which provides identical results for
 980 TCs 1-6), when $N_{MC} = 10,000$.

981

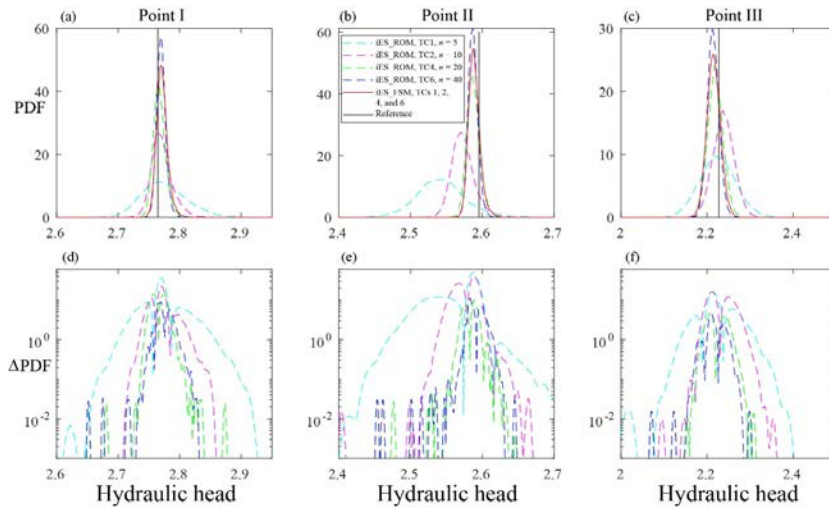


982

983 Fig. 7-8 Empirical PDFs of (a) x_{1,q_s} , (b) x_{2,q_s} , and (c) $\ln q_s$ at the final outer

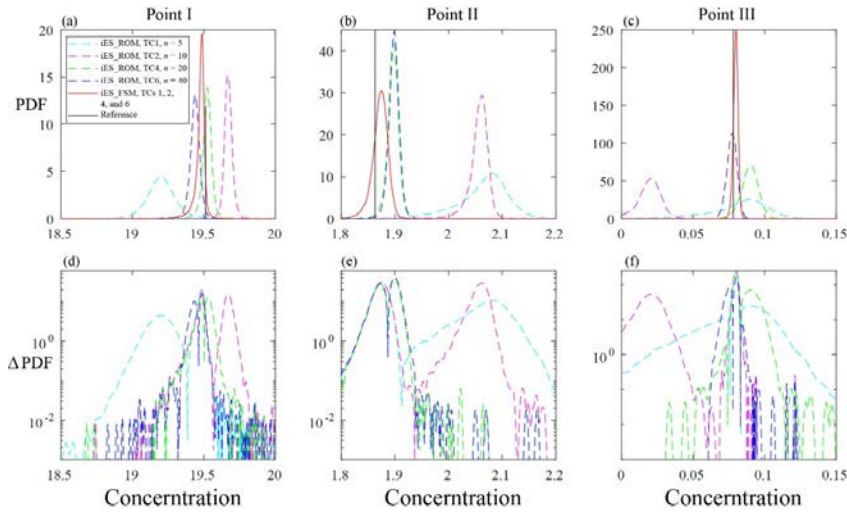
984 iteration through iES_ROM considering various values of n (equal to 5, 10, 20, and
 985 30 for TCs 1, 2, 4, and 6, respectively) and iES_FSM (which provides identical results
 986 for TCs 1, 2, 4, and 6) when $N_{MC} = 10,000$; corresponding reference values are
 987 indicated by black dashed lines.

988



989

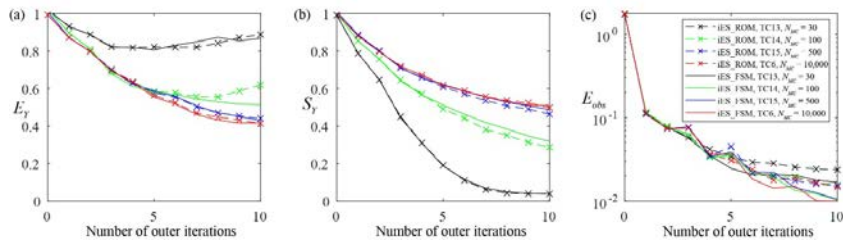
990 Fig. 8-9 Empirical PDFs of hydraulic head at points (a) I, (b) II, and (c) III (see Fig.
 991 4-2) at the final outer iteration obtained through iES_ROM considering different
 992 values of n (equal to 5 (cyan dashed curve), 10 (magenta), 20 (green), and 30 (blue)
 993 for TCs 1, 2, 4, and 6, respectively) and iES_FSM (red solid curve; identical results
 994 for TCs 1, 2, 4, and 6) when $N_{MC} = 10,000$ (corresponding reference values are
 995 indicated by black vertical lines); logarithmic absolute difference between PDFs
 996 obtained through iES_ROM and iES_FSM at points (a) I, (b) II, and (c) III.
 997



999

1000 Fig. 9-10 Empirical PDFs of solute concentration at points (a) I, (b) II, and (c) III (see
 1001 Fig. 4-2) at the final outer iteration obtained through iES_ROM considering various
 1002 values of n (equal to 5 (cyan dashed curve), 10 (magenta), 20 (green), and 30 (blue)
 1003 for TCs 1, 2, 4, and 6, respectively) and iES_FSM (red solid curve; results coincide
 1004 for TCs 1, 2, 4, and 6) when $N_{MC} = 10,000$ (corresponding reference values are
 1005 indicated by black lines); logarithmic absolute difference between PDFs obtained
 1006 through iES_ROM and iES_FSM at points (a) I, (b) II, and (c) III.

1007



1008

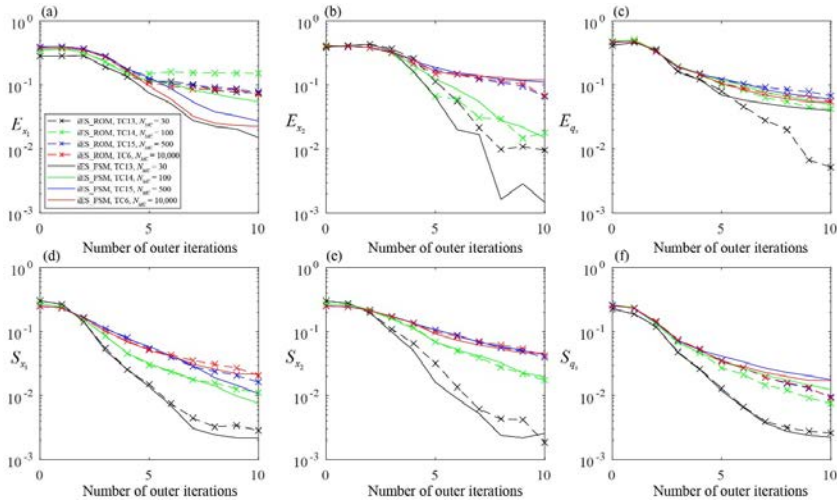
1009

1010

1011

1012

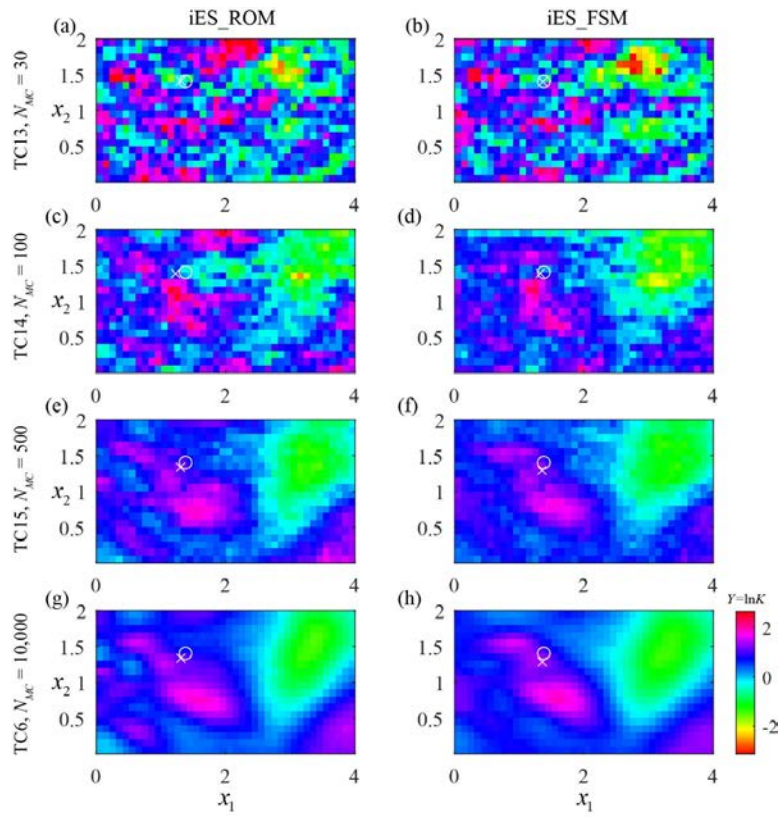
Fig. 40-11 Values of (a) E_Y , (b) S_Y , and (c) E_{obs} versus the number of outer iterations obtained through iES_ROM with $n = 30$ and iES_FSM considering various values of N_{MC} (equal to 30, 100, 500, and 10,000 for TCs 13-15 and 6, respectively).



1013

1014 Fig. 4-12 Values of (a) E_{x_1} , (b) E_{x_2} , (c) E_{q_1} , (d) S_{x_1} , (e) S_{x_2} , and (f) S_{q_1} versus

1015 the number of outer iterations obtained through iES_ROM with $n = 30$ and iES_FSM
 1016 considering various values of N_{MC} (equal to 30, 100, 500, and 10,000 for TCs 13-15
 1017 and 6, respectively).
 1018



1019

1020

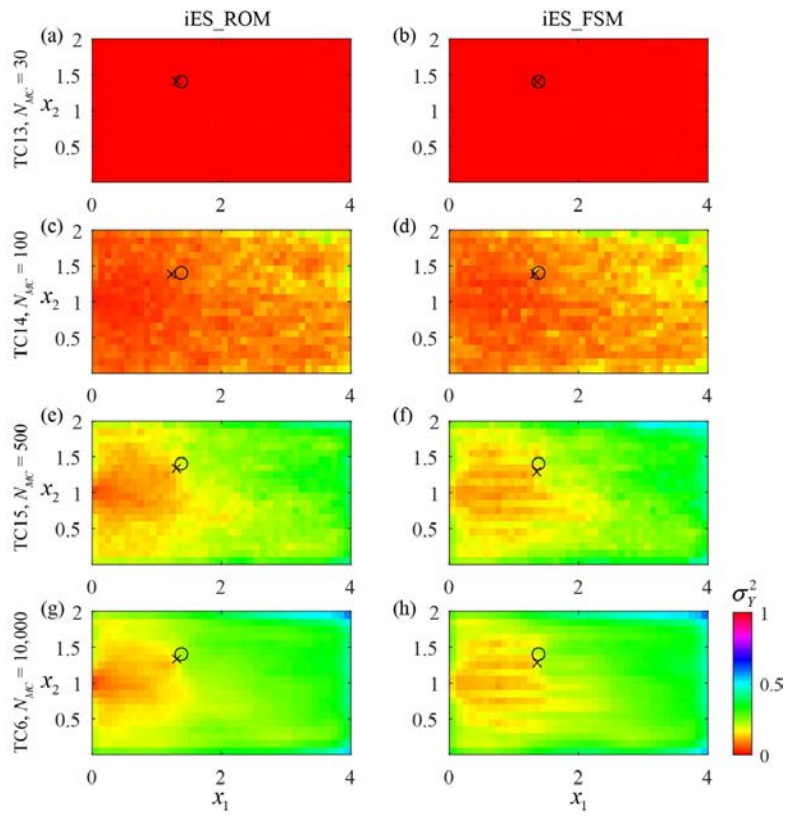
1021

1022

1023

1024

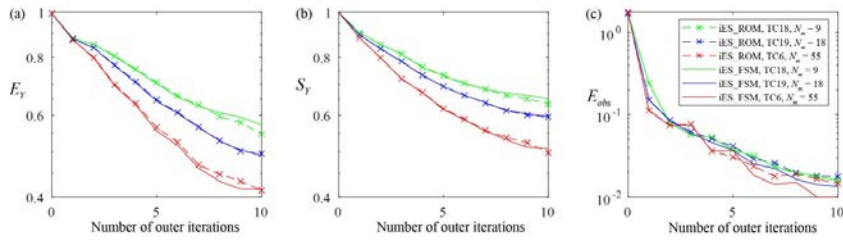
Fig. 42-13 Estimated (ensemble) mean Y fields at the final outer iteration obtained through iES_ROM with $n = 30$ (left column) and iES_FSM (right), considering $N_{MC} = 30$ (first row), 100 (second), 500 (third), and 10,000 (bottom) for TCs 13-15 and 6, respectively).



1025

1026 Fig. 13-14 Estimated (ensemble) Y variance fields at the final outer iteration obtained
 1027 through iES_ROM with $n = 30$ (left column) and iES_FSM (right), considering $N_{MC} =$
 1028 30 (first row), 100 (second), 500 (third), and 10,000 (bottom) (corresponding to TCs
 1029 13-15 and 6, respectively).

1030



1031

1032

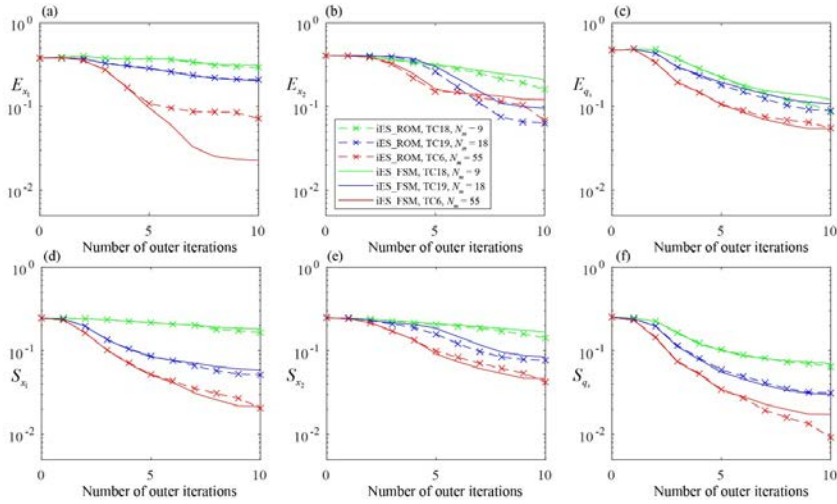
1033

1034

1035

1036

Fig. 14-15 Values of (a) E_Y , (b) S_Y , and (c) E_{obs} versus the number of outer iterations obtained through iES_ROM with $n = 30$ and iES_FSM considering $N_m = 9, 18,$ and 55 (corresponding to TCs 18, 19, and 6, respectively)



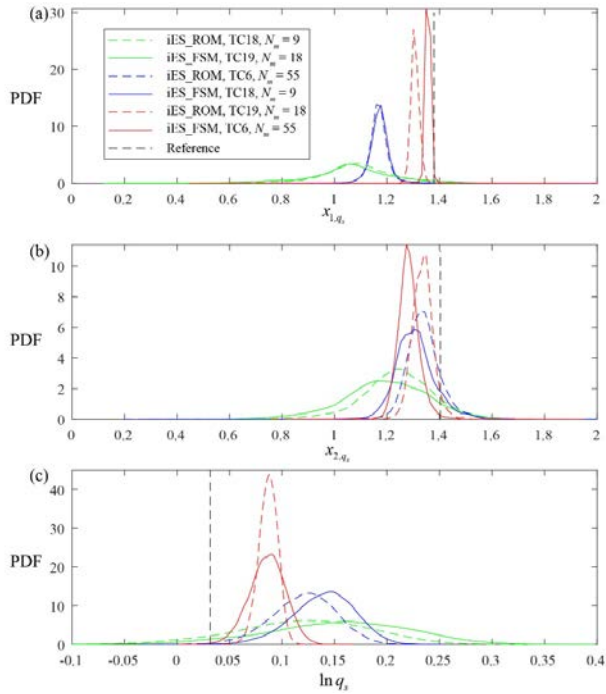
1037

1038 Fig. 45-16 Values of (a) E_{x_1} , (b) E_{x_2} , (c) E_{q_1} , (d) S_{x_1} , (e) S_{x_2} , and (f) S_{q_1} versus

1039 the number of outer iterations obtained through iES_ROM with $n = 30$ and iES_FSM,

1040 considering $N_m = 9, 18,$ and 55 (corresponding to TCs 18, 19, and 6, respectively).

1041



1042

1043

Fig. 16-17 Empirical PDFs of (a) x_{1,q_s} , (b) x_{2,q_s} , and (c) $\ln q_s$ at the final outer iteration through iES_ROM with $n = 30$ and iES_FSM (corresponding reference values indicated by black dashed lines) considering $N_m = 9, 18,$ and 55 (corresponding to TCs 18, 19, and 6, respectively).

1047

1048

Tables

1049 Table 1 Overview of the key settings of the test cases (TCs) analyzed. All TCs are

1050 characterized by a zero mean and unit variance of the Y reference field; μ and σ_Y^2 1051 denote the mean and variance of the initial ensemble of the Y fields, respectively.

Group A	Test Case	TC1, TC7	TC2, TC8	TC3, TC9	TC4, TC10	TC5, TC11	TC6, TC12
	n	5	10	15	20	25	30
	Known $q_s(\mathbf{x})$ or not	No, Yes	No, Yes	No, Yes	No, Yes	No, Yes	No, Yes
	Approach	iES_FSM and iES_ROM					
Group B	Test Case	TC13	TC14	TC15	TC6		
	N_{MC}	30	100	500	10,000		
	Approach	iES_FSM and iES_ROM					
Group C	Test Case	TC16	TC17	TC18	TC19	TC6	
	σ_{obs}	0.001	0.1	0.01	0.01	0.01	
	N_m	55	55	9	18	55	
	Approach	iES_FSM and iES_ROM					
Group D	Test Case	TC20	TC21	TC22	TC23	TC6	
	μ	-0.5	1.5	0.5	0.5	0.5	
	σ_Y^2	1.0	1.0	0.01	2.0	1.0	
	Approach	iES_FSM and iES_ROM					

	Test Case	TC24	TC25	TC26	TC27	TC28	TC6
Group E	N_{st}	30	100	300	500	1,000	10,000
	Approach	iES_ROM					

1052

1053

1054 Table 2 Values of E_Y , S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s} at the end of
 1055 the iteration procedure for TC16, TC6, and TC17 obtained through iES_ROM and
 1056 iES_FSM.

	TC16	TC6	TC17	TC16	TC6	TC17	TC16	TC6	TC17
	E_Y			S_Y			E_{obs}		
iES_ROM	0.41	0.41	0.53	0.50	0.50	0.62	0.01	0.02	0.07
iES_FSM	0.41	0.42	0.52	0.51	0.51	0.60	0.01	0.01	0.07
	E_{x_1}			E_{x_2}			E_{q_s}		
iES_ROM	0.07	0.07	0.15	0.08	0.07	0.04	0.05	0.06	0.13
iES_FSM	0.02	0.02	0.06	0.13	0.12	0.09	0.05	0.05	0.11
	S_{x_1}			S_{x_2}			S_{q_s}		
iES_ROM	0.02	0.02	0.04	0.04	0.04	0.08	0.01	0.01	0.03
iES_FSM	0.02	0.02	0.04	0.05	0.05	0.07	0.02	0.02	0.03

1057

1058 Table 3 Values of E_Y , S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s} at the end of
 1059 the iteration procedure for TC6, TC20, and TC21 obtained through iES_ROM and
 1060 iES_FSM.

Test Case	TC20	TC6	TC21	TC20	TC6	TC21	TC20	TC6	TC21
	E_Y			S_Y			E_{obs}		
iES_ROM	0.51	0.41	0.50	0.60	0.50	0.52	0.02	0.02	0.03
iES_FSM	0.44	0.42	0.52	0.55	0.51	0.56	0.01	0.01	0.03
	E_{x_1}			E_{x_2}			E_{q_s}		
iES_ROM	0.03	0.07	0.12	0.10	0.07	0.10	0.06	0.06	0.08
iES_FSM	0.01	0.02	0.12	0.10	0.12	0.17	0.05	0.05	0.11
	S_{x_1}			S_{x_2}			S_{q_s}		
iES_ROM	0.05	0.02	0.03	0.08	0.04	0.06	0.03	0.01	0.01
iES_FSM	0.03	0.02	0.04	0.06	0.05	0.06	0.02	0.02	0.02

1061

1062 Table 4 Values of E_Y , S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s} at the end of
 1063 the iteration procedure for TC6, TC22, and TC23 obtained through iES_ROM and
 1064 iES_FSM.

Test Case	TC22	TC6	TC23	TC22	TC6	TC23	TC22	TC6	TC23
	E_Y			S_Y			E_{obs}		
iES_ROM	0.47	0.41	0.46	0.06	0.50	0.77	0.01	0.02	0.02
iES_FSM	0.43	0.42	0.47	0.06	0.51	0.80	0.01	0.01	0.01
	E_{x_1}			E_{x_2}			E_{q_s}		
iES_ROM	0.07	0.07	0.09	0.15	0.07	0.12	0.02	0.06	0.08
iES_FSM	0.02	0.02	0.05	0.22	0.12	0.17	0.02	0.05	0.08
	S_{x_1}			S_{x_2}			S_{q_s}		
iES_ROM	0.003	0.02	0.05	0.004	0.04	0.09	0.003	0.01	0.02
iES_FSM	0.002	0.02	0.04	0.003	0.05	0.08	0.003	0.02	0.03

1065

1066

1067 Table 5 Percentage differences between the values of the selected metrics (i.e., E_Y ,
 1068 S_Y , E_{obs} , E_{x_1} , E_{x_2} , E_{q_s} , S_{x_1} , S_{x_2} , and S_{q_s}) at the end of the iteration procedure
 1069 for TCs 24-28 obtained through iES_ROM (values corresponding to TC6 are taken as
 1070 references).

Test Case	E_Y	S_Y	E_{obs}	E_{x_1}	S_{x_1}	E_{x_2}	S_{x_2}	E_{q_s}	S_{q_s}
TC24	11.88	6.34	21.60	44.17	58.50	25.04	34.75	32.29	55.20
TC25	10.16	3.66	6.76	27.83	35.54	13.24	9.58	25.01	37.58
TC26	7.40	1.97	13.00	22.10	16.89	35.54	11.26	13.06	15.09
TC27	4.58	0.14	2.66	11.42	0.17	22.93	1.22	17.74	3.37
TC28	0.50	0.21	4.18	17.19	1.33	31.33	5.86	14.03	0.07

1071

1072