1 **Efficient probabilistic hydrogeophysical prediction using maximum covariance analysis**

2 **based metamodel of the forward problem**

3

4

5 Erasmus K. Oware[1],

6 Jef Caers[2]

7 Thomas Hermans[3]

8

9 [1]*Department of Earth Sciences, SUNY at Buffalo, 126 Cooke Hall, Buffalo, NY 14260.*

10 [2]*Department of Geological Sciences, Stanford University, California, USA.*

11 [3]*Department of Geology, Ghent University, Belgium.*

12

13

14

15

16 **Key points**

17 • Improving computational efficiency of Bayesian inversion and sensitivity analysis

18   problems in hydrogeology and hydrogeophysics.

19 • Finding coupled patterns between two space-time domains using maximum covariance

20   analysis (MCA).

21 • First MCA-based metamodeling in hydrogeology and hydrogeophysics.

22

23

24

25

**Abstract**

Bayesian inverse modeling (BIM) provides a framework for uncertainty assessment in estimation of earth models from multiple sources of limited, noisy, and incomplete data. BIM, however, requires repeated solutions of computationally expensive full-physics forward (FPF) runs, rendering BIM overwhelmingly computationally prohibitive and intractable in high-dimensional problems. We present a novel maximum covariance analysis (MCA)-based metamodel to approximate the FPF problem by capturing coupled patterns from small number of mutual observations between the parameter and data fields. We implement the strategy as a difference MCA prediction-focused approach (DMCA-PFA) and test it on electrical resistivity tomography (ERT) field data acquired during a heat-tracer experiment. We also invert the data using the conventional smoothness-constrained inversion (SCI). The DMCA-PFA outperformed SCI in estimation of realistic plume morphologies and direct temperature validations. We show an excellent MCA-based approximations of the FPF simulated resistances, with the potential to reduce CPU-time of BIM by over 99%.

**Plain Language Summary**

Mathematical modeling is essential for efficient management of groundwater and energy resources. Such mathematical procedures require input parameters of some earth properties. Since the earth cannot be sampled exhaustively to obtain the input parameters at all locations, inverse modeling is performed to derive the input parameters at unsampled locations from observations at limited sampled locations. The inverse modeling procedure also requires forward modeling, which involves mathematical modeling of a process response given the input earth properties. The severally repeated forward modeling task, however, can become overwhelming computationally expensive, especially when the modeling involves several spatially distributed

parameters, which limits our inverse modeling capabilities. We present a novel maximum

covariance analysis (MCA)-based strategy to approximate the forward problem in order to

circumvent the tremendous computational costs of computing several forward solutions. We test

the strategy on electrical resistivity tomography (ERT) field data acquired during a heat-tracer

experiment. We also invert the data using the conventional smoothness-constrained inversion

(SCI). Our strategy outperformed SCI in estimation of realistic plume morphologies and direct

temperature validations. We show an excellent MCA-based approximations of the forward

problem, with the potential to reduce CPU-time of inverse modeling by over 99%.

## 1. Introduction

Inverse modeling in hydrogeology [e.g., *Zhou et al.,* 2014] and hydrogeophysics [e.g., *Binley et al.*, 2015] involves estimating some earth physical properties and processes from observational data. Such observations are typically noisy and limited in number and resolution coupled with incomplete understanding of coupled processes occurring at multiple spatio-temporal scales. Bayesian inversion [*Tarantola*, 2005] provides a framework to infer earth models from multiple sources of information, including noisy and incomplete data while allowing incorporation of prior information. It also allows uncertainty quantification, thereby enabling comprehensive interrogation of the estimates. Due to the usually high-dimensional and complex non-linear models, the solution necessitates Markov chain Monte Carlo (McMC) sampling [*Hansen et al.*, 2016]. Bayesian McMC, however, requires repeated solutions of computationally expensive full-physics forward runs, which renders them computationally prohibitive and intractable in high-dimensional problems.

There is growing interest in the use of proxy-models of the forward solution in applications that demand several computationally expensive forward runs, such as Bayesian McMC and

74    sensitivity analysis. There are two broad categories of proxy-models, the lower-fidelity models

75    (LF-models) and metamodels [*Linde et al*., 2017]. The LF-models are usually physics-based and

76    rely on model simplifications to speed-up the forward runs. The simplification can be achieved

77    by approximating some aspects of the physics or by ignoring them completely [e.g., *Josset et al.*,

78    2015], or by performing the forward runs on a coarser grid [*Arridge et al*., 2015]. While LF-

79    models speed-up the forward runs, they are less accurate. The metamodels, in contrast, are data-

80    driven proxies involving Monte Carlo simulations of relatively small number of ensembles

81    obtained from the full-physics (high-fidelity) forward simulations. There are several methods for

82    developing the proxy-models from the small number of high-fidelity ensembles [e.g., *Khu and*

83    *Werner*, 2003; *Myers et al*., 2016]. In the context of metamodeling in hydrogeology, the

84    polynomial chaos expansion (PCE) [*Beck et al*., 2014] and reduced-order models (ROM) [e.g.,

85    *Liu et al.*, 2013] are commonly applied. In the PCE metamodeling, polynomial approximations

86    of the forward problem are constructed over the support of the prior distributions [e.g., *Marzouk*

87    *and Xiu*, 2009]. PCE, however, suffers in high-dimensional parameter spaces since the number

88    of PCE terms increases dramatically with increased number of input parameters. PCE also

89    underperforms when the input random field is highly heterogeneous. The ROM approach

90    constructs orthogonal bases from the small number of high-fidelity ensembles (snapshots) and

91    then employs the orthogonal bases as projection matrices to map the high-dimensional target

92    system into a low-dimensional subspace. Unlike the current ROM techniques that consider

93    orthogonal bases of a single state space, such as hydraulic heads [.e.g., *Liu et al*., 2013] or solute

94    concentrations [e.g., *Oware et al.*, 2013; *Oware and Moysey*, 2014], our metamodeling relies on

95    coupled (joint) patterns between the model parameter and data domains, making the strategy

96    particularly well suited for metamodeling in model-data integration and sensitivity analysis

4

97   applications. *Linde et al.* [2017] provides an excellent review of proxy-modeling in

98   hydrogeology and hydrogeophysics.

99      Furthermore, to avoid the overwhelmingly computationally expensive Bayesian inversion,

100  *Scheidt et al.* [2015] proposed the prediction-focused approach (PFA) to predict hydrologic

101  variables directly from geophysical measurements without the full inversion of the data nor post-

102  inversion petrophysical transform. The PFA uses surrogate models to derive a linearized,

103  statistical relationship between the data and the target parameters [e.g., *Hermans et al.*, 2016b].

104  While PFA shows a lot of promise, it predicts the target hydrologic models directly from the

105  geophysical data without iteratively fitting the data, which might limit its robustness to

106  reconstruct complex hydrologic features that are not well represented in the prior samples

107  [*Oware et al.*, in-press]. The current PFA framework [e.g., *Satija and Caers*, 2015], moreover,

108  uses canonical correlation analysis (CCA) to capture the coupled relationship between the two

109  (model-data) space-time domains. The use of multivariate statistical tools, such as maximum

110  covariance analysis (MCA) to capture coupled patterns between two space-time parameter fields

111  for the purpose of forecasting is widely used in climate science [e.g., *von Storch and Zwiers*,

112  1999]. MCA is simply singular value decomposition (SVD) of the cross-covariance between the

113  two space-time domains. *Bretherton et al.* [1992], for instance, found CCA to be uncompetitive

114  compared to MCA due to high sampling variability unless the coupled signal was highly

115  localized.

116     We propose here a difference MCA-PFA (DMCA-PFA) scheme to advance the PFA

117  framework with two key contributions: 1) implement PFA in the MCA coupled space by actually

118  fitting the data in a Bayesian sense, and 2) develop an MCA-based metamodel of the geophysical

119  forward problem. We also outline a strategy to calibrate and account for metamodel-discrepancy

120    in the proxy-approximation. We illustrate the performance of the DMCA-PFA on a field-scale

121    geoelectrical data acquired during a heat-tracer experiment. We intend to show that the MCA-

122    based metamodel presents a key contribution toward improving the computational efficiency of

123    stochastic inversion and sensitivity analysis procedures in hydrogeology and hydrogeophysics.

124

125    **2. Methods**

126    2.1 Overview of Maximum Covariance Regression

127    Consider two parameter fields, a geophysical data field, $\mathbf{d} \in \mathbf{R}^q$ and a hydrologic model

128    parameter space, $\mathbf{h} \in \mathbf{R}^p$, where $q$ and $p$ are the number of geophysical data points and

129    hydrologic model parameters, respectively. A linear multivariate regression between the two

130    domains can be expressed as: $\qquad \mathbf{d} = \mathbf{Ah},$ $\qquad\qquad\qquad$ (1)

131    where $\mathbf{A}$ is a regression matrix. The mapping in Equation 1 usually involves complex, non-linear

132    relationships including spatially dependent petrophysical transformation between the $\mathbf{d}$ and $\mathbf{h}$. To

133    linearize such complex relationships, we propose to use maximum covariance analysis (MCA).

134    A good treatment of MCA is provided by *von Storch and Zwiers* [1999]. To accomplish this, we

135    construct the data matrix $\mathbf{D} \in \mathbf{R}^{q \times n}$ and the model parameter matrix $\mathbf{H} \in \mathbf{R}^{p \times n}$ from Monte

136    Carlo simulations of $n$ number of mutual observations (snapshots) between $\mathbf{d}$ and $\mathbf{h}$. If $\mathbf{HH}^T$ is

137    invertible, then $\mathbf{A}$ in Equation 1 can be factorized in terms of MCA projections [e.g., *Tippett et*

138    *al.*, 2008]: $\qquad\qquad\qquad \mathbf{A} = \mathbf{U\Lambda V}^T(\mathbf{HH}^T)^{-1},$ $\qquad\qquad\qquad$ (2)

139    where $T$ denotes transpose, $\mathbf{U\Lambda V}^T$ is the SVD of $\mathbf{DH}^T$, $\mathbf{\Lambda} \in \mathbf{R}^{p \times p}$ is a diagonal matrix of singular

140    values, and $\mathbf{U} \in \mathbf{R}^{q \times p}$ and $\mathbf{V} \in \mathbf{R}^{p \times p}$ are the left (data) and right (hydrologic model) coupled

141    patterns, respectively. From Equations 1 and 2, we recast $\mathbf{d}$ as:

142    $\qquad\qquad\qquad\qquad \mathbf{d} = \mathbf{U\Lambda V}^T(\mathbf{HH}^T)^{-1}\mathbf{h} + \boldsymbol{\varepsilon},$ $\qquad\qquad\qquad$ (3)

143   where $\boldsymbol{\varepsilon}$ is the metamodel-discrepancy that accounts for the inexactness of the MCA-based

144   approximation of the high-fidelity $\mathbf{d}$.  A consequence of Equation 3 is that if we learn the mutual

145   behavior between $\mathbf{d}$ and $\mathbf{h}$ and the metamodel-discrepancy structure from training surrogates $\mathbf{D}$

146   and $\mathbf{H}$,  then we can directly predict $\mathbf{d}$ associated with any given $\mathbf{h}$ without the need for the

147   typically computationally expensive geophysical forward simulations. We only need to run the

148   high-fidelity forward simulations only $n$ number of times.

149       In the event that $\mathbf{HH}^T$ is not invertible directly, an inverted version can be approximated via

150   SVD [e.g., *Castleman*, 1996]:

$$(\mathbf{HH}^T)^{-1} = \left(\mathbf{U}_h \boldsymbol{\Lambda}_h \mathbf{V}_h{}^T\right)^{-1} \approx \mathbf{V}_h \boldsymbol{\Lambda}_h{}^{-1} \mathbf{U}_h{}^T \tag{4}$$

152   where $\mathbf{U}_h \boldsymbol{\Lambda}_h \mathbf{V}_h{}^T$ is SVD of $\mathbf{HH}^T$, $\boldsymbol{\Lambda}_h{}^{-1}$ is a diagonal matrix with its diagonal elements equal to

153   $1/\Lambda_{ij}$ and $\Lambda_{ij}$ are the diagonal elements of $\boldsymbol{\Lambda}_h$.

154

155   2.2 Difference Maximum Covariance Analysis Prediction-Focused Approach

156       Difference inversion [*LeBrecque and Yang*, 2001] has become increasingly appealing for

157   geophysical monitoring of hydrogeological processes because inverting on the background

158   differenced data results in rapid convergence, ability to detect small changes, eliminate

159   systematic errors, and reduce inversion artifacts. Hence, we test the strategy as a difference

160   maximum covariance analysis prediction-focused approach (DMCA-PFA). We apply Bayes'

161   rule for the problem of estimating the posterior distribution of $\mathbf{h}$ from observed data, $\mathbf{d}_{obs}$.

162   Specifically,       $$\mathbf{h}_{post} = \mathbf{h}_{prior} L(\mathbf{h}|\mathbf{d}_{obs}) \tag{5}$$

163   where $\mathbf{h}_{post}$ and $\mathbf{h}_{prior}$ are the posterior and prior distributions of $\mathbf{h}$, respectively, and $L(\cdot)$ is the

164   likelihood, which evaluates the probability of a proposed $\mathbf{h}$ given $\mathbf{d}_{obs}$. We compute the

165   likelihood as a multivariate Gaussian error distribution, i.e.,

$$L(\mathbf{h}|\mathbf{d}_{obs}, \mathbf{W}_d) = \exp\left[-\frac{1}{2}(\mathbf{e}^T * \mathbf{W}_d * \mathbf{e})\right], \tag{6}$$

167 where $\mathbf{W}_d$ is the data weight matrix and $\mathbf{e}$ is the data misfit. To implement Equation 6 in a

168 difference inversion framework, we express $\mathbf{e}$ as:

$$\mathbf{e} = [\mathbf{d}_t - \mathbf{d}_0] - [f(\mathbf{h}_t) - f(\mathbf{h}_0)], \tag{7}$$

170 where $\mathbf{d}_t$ and $\mathbf{d}_0$ represent the data at the time-step of interest and background, respectively. The

171 terms $f(\mathbf{h}_t)$ and $f(\mathbf{h}_0)$ are, respectively, the predicted data associated with a proposed model

172 and the model obtained from the classical inversion of the background data. To advance the PFA

173 framework by actually fitting the observed data in a computationally efficient manner, we

174 circumvent the high-fidelity geophysical forward runs in Equation 7 by directly predicting the

175 data for any given $\mathbf{h}$ according to Equation 3. Hence,

$$\mathbf{e} = [\mathbf{d}_t - \mathbf{d}_0] - [\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T(\mathbf{HH}^T)^{-1}\mathbf{h}_t - \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T(\mathbf{HH}^T)^{-1}\mathbf{h}_0 + (\boldsymbol{\varepsilon}_t - \boldsymbol{\varepsilon}_0)]. \tag{8}$$

177 Similar to *Hermans et al.* [2016b], Equation 8 predicts the hydrologic model directly from the

178 geophysical data, thereby avoiding post-inversion petrophysical transformation. There is

179 growing popularity in estimation algorithms that proceed in the reduced-dimensional space due

180 to their computational stability and efficiency [e.g., *Banks et al.*, 2000]. Hence, we express the

181 target parameters, $\mathbf{h}_t$, as a linear combination of its basis vectors, $\mathbf{B}$, and expansion coefficients,

182 $\mathbf{c}$, i.e., $\mathbf{h}_t = \mathbf{B}\mathbf{c}_t$ [e.g., *Oware et al.*, 2013]. Therefore,

$$\mathbf{e} = [\mathbf{d}_t - \mathbf{d}_0] - \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T(\mathbf{HH}^T)^{-1}[\mathbf{B}\mathbf{c}_t - \mathbf{h}_0] - [\boldsymbol{\varepsilon}_t - \boldsymbol{\varepsilon}_0]. \tag{9}$$

184 It should be emphasized that the basis, $\mathbf{B}$, is constructed from SVD of $\mathbf{H}$ and, therefore, $\mathbf{B}$ is not

185 the same as $\mathbf{V}$. While $\mathbf{V}$ captures the MCA coupled patterns between $\mathbf{H}$ and $\mathbf{D}$, $\mathbf{B}$ represents

186 orthogonal bases of $\mathbf{H}$ only.

187    To obtain prior distributions for the metamodel-discrepancy, $\boldsymbol{\varepsilon}_t$, in Equation 9, we perform

188 MCA-based approximations (Equation 3) of the geophysical data associated with the *n* number

189 of training samples (**H**), to construct the data matrix $\mathbf{D}_{mca} \in \mathbf{R}^{q \times n}$. We then compute the prior

190 error, $\boldsymbol{\varepsilon}_{prior} \in \mathbf{R}^{q \times n}$, as the discrepancies between the predicted data of the high-fidelity

191 geophysical forward runs and those of the MCA-based approximations, i.e.,

192 
$$\boldsymbol{\varepsilon}_{prior} = \mathbf{D} - \mathbf{D}_{mca}. \tag{10}$$

193 This implies that there are $n$ realizations of the approximation errors for each geophysical data

194 point, which defines prior distributions for sampling the metamodel-discrepancy, $\boldsymbol{\varepsilon}_t$, for each

195 geophysical data point. We compute $\boldsymbol{\varepsilon}_0$ in Equation 9 as the residuals between the predicted data

196 of the high-fidelity geophysical forward run and that of the MCA-based approximation obtained

197 from the inverted background geophysical model.

198

199 2.3 Complete Overview of DMCA-PFA

200 We now present the full workflow of the sampling of the posterior distribution of **h** (Figure

201 1):

202 1) Perform Monte Carlo simulations to generate training set (TS) of multiple realizations of the

203 physics of the target hydrologic process (e.g., captures outcomes of multiple rates of advection

204 and multiple scales of dispersion and plume morphologies). Collect all simulated time-lapse

205 models into a single library of hydrologic (space-time) TS (**H**).

206 2) To obtain mutual observations between **h** and **d**, perform petrophysical transformation of each

207 **h** into geophysical properties and run geophysical forward simulations to predict **d** associated

208 with each **h**. Collect all simulated **d** into a matrix of geophysical TS, **D**. Now, **D** and **H** constitute

209 mutual observations between the two parameter fields, **h** and **d**.

210 3) Perform MCA of **D** and **H** to obtain coupled patterns between the two fields. Also, SVD of **H**

211 produces orthogonal bases, **B**.

212   4) To obtain prior distributions for sampling the expansion coefficients, **c**, project **H** onto **B**, i.e.,

213   $\mathbf{c}_{prior} = \mathbf{B}^T\mathbf{H}$. Obtaining the prior coefficients from the physics-based TS imposes physics-

214   based parameter bounds for the coefficients in an attempt to produce physically realistic plume

215   morphologies [*Oware et al*., 2018].

216   5) Propose coefficients from **c***prior*. We accept or reject the proposed coefficients based on the

217   classical Metropolis-Hastings acceptance rule [*Metropolis et al*., 1953; *Hastings*, 1970]. The

218   posterior coefficients are then mapped onto **B** to obtain multiple realizations of the target, i.e.,

219   $\mathbf{h}_{post} = \mathbf{B}\mathbf{c}_{post}$. Note, Step 5 is simply the standard McMC sampling parameterized in the

220   reduced-dimensional space. It also uses the MCA-based metamodel without performing the

221   typically computationally expensive geophysical forward runs.

222

223   **3. Application to Field Data**

224     We demonstrate the efficacy and efficiency of the DMCA-PFA algorithm on a field-scale

225   heat-tracer experiment conducted in an alluvial aquifer and monitored with cross-well ERT

226   (XBh-ERT). Details of the heat-tracer and XBh-ERT surveys are outlined in *Hermans et al*.

227   [2015]. To summarize, water was continuously pumped to induce GW flow toward the pumping

228   well. Hot water was then injected continuously in an injection well for 24 hours. Changes in

229   electrical conductivity were monitored in a XBh-ERT panel across the GW flow direction. We

230   invert the first six time-lapse resistances acquired at 6 h, 12 h, 18 h, 21.5 h, 25 h, and 30 h after

231   the commencement of the heat injection. Each time-step inversion involves only 410

232   quadrupoles. Direct temperatures were also monitored in two piezometers, pz14 and pz15, for

233   validation of the ERT predicted temperatures.

234    We first performed Monte Carlo simulations to obtain a training set (TS) tuned to the physics

235    of the presupposed heat-tracer test. We used the same 3,000 (500 hydrologic models x 6 time-

236    steps) temperature TS (**H**) employed by *Hermans at al*. [2018]. The TS was obtained via Monte

237    Carlo simulations of the heat tracing experiment for 500 different GW models, considering

238    uncertainties in the underlying hydrogeologic and transport properties. Through petrophysical

239    transformations, we converted each temperature distribution, **h**, into resistivity models and ran

240    resistivity forward simulations to obtain resistance data, **d**, associated with each **h**. A collection

241    of all **d** comprise the geophysical training data (**D**). We then performed MCA of **D** and **H** to

242    construct the coupled patterns between the two fields.  Figure 2 shows the first 5 dominant MCA

243    coupled patterns between the hydrological (log(**h**)) and  geophysical data (**d**)  spaces constructed

244    from the 3,000 mutual observations between the two fields. Figure 2, essentially, depicts how the

245    two fields covary such that given any resistance measurements, **d**, we should be able to leverage

246    the prior coupled behavior to predict its associated temperature distributions, **h**.  For comparison,

247    we also inverted all the datasets using the classical smoothness-constrained inversion (SCI). We

248    applied the 2.5D ERT inversion code CRTomo [*Kemna*, 2000] for all resistivity forward

249    simulations and the SCI. We utilized the petrophysical relationship presented by *Hermans et al*.

250    [2015] and the parameters presented therein for all conversions of ERT into thermograms.

251

## 4. Results and Discussion

252    

*4.1 MCA-based metamodeling*

253    

254    For the inversion of $\mathbf{HH}^T$ in Equation 3, we used Equation 4 since $\mathbf{HH}^T$ was not invertible in

255    the case study presented here. Histogram analyses (not shown) of the metamodel-discrepancies

256    of the individual data points (Equation 10) reveal that the errors are not normally distributed.

257    Hence, we assumed no knowledge of the prior error distributions and sampled uniformly over

258    the interval of the prior errors of each data point for $\varepsilon$ in Equation 3. To assess the performance

259    of the MCA-based metamodel, we applied the high-fidelity resistivity forward simulation and

260    MCA-based approximation to estimate resistances from resistivity tomograms obtained from

261    smoothness-constrained inverison of the observed resistance data at three time steps, 12 hours

262    (t2), 21.5 hours (t4), and 30 hours (t6).

263    Figure 3 shows the scatter plots of the full resistivity forward simulated resistances against

264    those of the MCA-based metamodel for the three time steps. The coefficients of determination

265    ($R^2$) for the MCA-based metamodel for t2, t4, and t6 are 0.9996, 0.9990, and 0.9988,

266    respectively. The $R^2$s indicate marginal deterioration of the MCA metamodel with increasing

267    time-steps. Nevertheless, there is almost a perfect one-to-one MCA proxy-approximations of all

268    the high-fidelity forward simulated resistances, indicating high approximation accuracy of the

269    MCA metamodeling in the examples considered here. It takes ~2.64 seconds of CPU-time to

270    complete each high-fidelity resistivity forward simulation. This implies that Bayesian inversion

271    involving about 300,000 iterations, for instance, will require ~13,200 minutes of CPU-time. The

272    Bayesian inversion with the MCA-based metamodel (DMCA-PFA) presented in the next section

273    takes ~35 minutes to complete 300,000 iterations. This represents a significant reduction in the

274    computional time of ~99%, considering the 3,000 high-fidelity forward runs needed to calibrate

275    the MCA-based metamodel and the ~35 minutes needed to complete the inversion.

276

277    *4.2 Posterior Prediction*

278    We ran the algorithm for 300,000 iterations for all of the six time-lapse profiles. We applied

279    20 bases, **B**, to reconstruct the 1092 full-dimensional space, resulting in over 98% truncation in

280    the dimensionality of the problem. The 20 selected basis vectors represented 99.8% of the total

281    variance in the TIs. The difference thermograms recovered for the 12h (t2), 21.5h (t4), and 30h

282    (t6) time-steps based on the SCI and the DMCA-PFA are presented in Figure 3. While both

283    strategies captured similar evolution (locations and spatial extents) of the heat plume (Figure 3

284    Columns 1-4), smoothing of the heat plume is less severe in our approach in contrast to

285    smoothing in the SC thermograms. The estimation of physically realistic plume morphologies

286    without excessive smoothing in our approach (Figure 3 Columns 2-4) is attributable to the use of

287    physics-based prior constraints as compared to the use of a generalized smoothness spatial filter

288    in the SCI. The standard deviation panels (Figure 3 Column 5) reveal the spatial variabilities of

289    uncertainty in the estimates. While uncertainty is expected to be low near the borehole locations

290    (extreme vertical ends) due to high cross-borehole resistivity data sensitivity near the ERT wells,

291    there appears to be generally high uncertainty in the recovered temperatures around 8-9 m

292    depths, especially near the left borehole corresponding to high amplitudes. This trend in the

293    estimated uncertainty reveal increasing difficulty of the algorithm to estimate high temperature

294    deviations from the background values. The ability of our strategy to accurately capture the

295    migration of the heat plume (different locations and morphologies) using the same set of basis-

296    constraints demonstrates the flexibility of the algorithm to recombine the bases in a manner that

297    honors each time-step ERT measurements.

298

299    *4.3 Model Validation with Direct Temperature Measurements*

300        Figure 5 outlines the validation of estimated temperature breakthrough curves at the two

301    piezometers, pz14 and pz15, respectively, located at (1.125 m, 9 m) and (2.25 m, 8.5 m) from the

302    left borehole.  The temporal behavior of the validation breakthrough curves were accurately

303    captured by both methods. Comparisons of the estimated temperatures with the direct

304    temperatures, however, indicate that DMCA-PFA outperformed SCI on almost all the direct

305    temperature measurements. The 90% confidence interval (CI) of the DMCA-PFA estimates

306    captured all the true temperatures with the true values seemingly well centered within the range

307    of the 90% CI. *Hermans et al*. [2018] concluded that a change of $1^{o}C$ produced only 2% change

308    in electrical conductivity for the data presented here. Such small changes are undetectable in

309    deterministic inversions [e.g., *Doetsch et al.*, 2012). *Hermans et al*. [2015] estimated the limit of

310    detection of ERT of this experiment at ~1.5 $^{o}C$ given the estimated noise level. Accounting for

311    the physics of the target process seems to improve the limit of detection in the DMCA-PFA.

312    Particularly, the direct temperature measurements at both pz14 and pz15 (Figures 5A and 5B)

313    show that 6 hours (t1) of heat injection resulted in a change of ~0.5 $^{o}C$, which was undetected by

314    the SCI since it is well below the ~1.5 $^{o}C$ ERT detection limit. Our approach, in contrast,

315    accurately estimated the small temperature change and captured the true values within 90% CI.

316

317    **5. Conclusion**

318        Inverse modeling is the foremost strategy for inferring earth properties and processes from

319    observational data in hydrogeology and hydrogeophysics. In spite of the numerous benefits of

320    stochastic inversion, deterministic inverse methods remain widely used due to their simplicity

321    and computational efficiency. We propose here a novel maximum covariance analysis (MCA)-

322    based metamodel to reduce the overwhelming computational costs of repeatedly computing the

323    full-physics of the forward problem in stochastic inversions. We construct the MCA-based

324    metamodel from coupled patterns captured from small number of ensembles of the joint

325    evolution of the parameter and data fields. Hence, the strategy accounts for the physics of the

326    target system on two fronts, the physics of the parameter and data acquisition systems. We

327    conclude that incorporating the physics of the target process improves estimation and produces

328     physically realistic target plumes without excessive smoothing in contrast to results obtained

329     from the conventional smoothness-constrained inversion. We found an excellent MCA-based

330     proxy-approximations of the full-physics forward simulated data, with the potential to reduce

331     CPU-time of Bayesian inverse procedures by over 99%. The MCA-based metamodel presents a

332     promising general framework to speed-up the computational efforts of hydrogeological and

333     geophysical applications that necessitate repeated computations of the full-physics of the forward

334     problem, such as high-dimensional Bayesian inversion and sensitivity analysis problems.

335

339

340     **References**

341     Arridge, S., J. Kaipio, V. Kolehmainen, M. Schweiger, E. Somersalo, T. Tarvainen, M.

342             Vauhkonen (2006), Approximation errors and model reduction with an application in

343             optical diffusion tomography. *Inverse Probblems,* 22 (1), 175−195.

344             http://dx.doi.org/10.1088/0266-5611/22/1/010.

345     Banks, H.T., M.L. Joyner, B. Winchesky, W.P. Winfree (2000), Nondestructive evaluation using

346             a reduced order computational methodology, *Inverse Problems,* 16, 1–17.

347     Beck, J., F. Nobile, L. Tamellini, R. Tempone (2014), A quasi-optimal sparse grids procedure

348             for groundwater flows. Spectral and High Order Methods for Partial Differential Equations -

349             ICOSAHOM 2012. *Lecture Notes in Computational Science and Engineering 95*. Springer,

350             pp. 1−16. http://dx.doi.org/10.1007/978-3-319-01601-6_1.

Bretherton, C.S., C. Smith, and J.M. Wallace (1992), An intercomparison of methods for finding
coupled patterns in climate data, *J. Climate*, 5, 541-560.

Binley, A., S. S. Hubbard, J. A. Huisman, A. Revil, D. A. Robinson, K. Singha, and L. D. Slater
(2015), The emergence of hydrogeophysics for improved understanding of subsurface
processes over multiple scales, *Water Resour. Res.*, 51, doi:10.1002/2015WR017016.

Castleman, K.R. (1996), Digital Image Processing: Prentice Hall, Inc., Upper Saddle River, NJ.

Doetsch, J., Linde, N., Vogt, T., Binley, A., & Green, A. G. (2012), Imaging and quantifying
salt-tracer transport in a riparian groundwater system by means of 3D ERT monitoring,
*Geophysics*, 77, B207‑B218. https://doi.org/10.1190/geo2012-0046.1

Hansen, T.M., K.S. Cordua, A. Zunino, K. Mosegaard, (2016), Probabilistic integration of
geo-information. In: Moorkamp, N.L.M., Leliévre, P.G., Khan, A. (Eds.), Integrated
Imaging of the Earth: Theory and Applications, John Wiley & Sons, Inc, pp. 93–116.

Hastings, W. (1970), Monte Carlo sampling methods using Markov chains and their
applications, *Biometrika*, 57, no. 1, 97.

Hermans, T, E. K. Oware, and J.K. Caers, 2016b, Direct prediction of spatially and temporally
varying physical properties from time-lapse electrical resistance data: Water Resources
Research, 52, no. 9, 7262-7283, doi.org/10.1002/2016WR019126.

Hermans, T., F. Nguyen, M. Klepikova, A. Dassargues, & J. Caers (2018), Uncertainty
quantification of medium-term heat storage from short-term geophysical experiments
using Bayesian evidential learning, *Water Resources Research*, 54.

Hermans, T., S. Wildemeersch, P. Jamin, P. Orban, S. Brouyere, A. Dassargues, and F.
Nguyen (2015), Quantitative temperature monitoring of a heat tracing experiment using
cross borehole ERT: *Geothermics*, 53, 14–26, doi:10.1016/j.geothermics.2014.03.013.

374　Josset, L., V. Demyanov, A.H. Elsheikh, I. Lunati (2015), Accelerating Monte Carlo Markov

375　　　　chains with proxy and error models, *Computer Geoscience,* 85, 38–48.

376　Kemna, A. (2000), Tomographic inversion of complex resistivity: theory and application: PhD

377　　　　Thesis, Bochum Ruhr University, Germany.

378　Khu, S.T., and M. Werner (2003), Reduction of Monte-Carlo simulation runs for uncertainty

379　　　　estimation in hydrological modelling, *Hydrol. Earth Syst. Sci.* 7 (5), 680–692.

380　LaBrecque, D.J., and X. Yang (2001), Difference inversion of ERT data: a fast inversion method

381　　　　for 3-D in-situ monitoring: Journal of Environmental and Engineering Geophysics, 6, no.

382　　　　2, 83 – 89.

383　Linde, N, D. Ginsbourgerb, J. Irving, F. Nobiled, and A. Doucet (2017), On uncertainty

384　　　　quantification in hydrogeology and hydrogeophysics, *Advances in Water Resources,* 110,

385　　　　166–181: http://dx.doi.org/10.1016/j.advwatres.2017.10.014.

386　Liu, X., Q. Zhou, J. Birkholzer, W.A. Illman (2013), Geostatistical reduced-order models in

387　　　　underdetermined inverse problems, *Water Resources Research*, 49 (10), 6587–6600.

388　Marzouk, Y., and D. Xiu (2009), A stochastic collocation approach to Bayesian inference in

389　　　　inverse problems. *Commun. Comput. Phys.* 6, 826–847. http://dx.doi.org/10.4208/

390　　　　cicp.2009.v6.p826.

391　Metropolis, N., A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller (1953), Equation of state

392　　　　calculations by fast computing machines, Journal of Chemical Physics, 21, 1087–1092.

393　Myers, R., D. Montgomery, C. Anderson-Cook (2016), Response Surface Methodology: Process

394　　　　and Product Optimization Using Designed Experiments. Wiley.

395 Oware, E. K., and S. M. J. Moysey (2014), Geophysical evaluation of solute plume spatial

396          moments using an adaptive POD algorithm for electrical resistivity imaging: Journal of

397          Hydrology, 517, 471-480. http://dx.doi.org/10.1016/j.jhydrol.2014.05.054

398 Oware, E. K., S. M. J. Moysey, and T. Khan (2013), Physically based regularization of

399          hydrogeophysical inverse problems for improved imaging of process-driven systems,

400          *Water Resources Research*, 49(10), 6238-6247. http://dx.doi.org/10.1002/wrcr.20462.

401 Oware, E. K., M. Awatey, T. Hermans and J. Irving (2018), Basis-Constrained Bayesian-McMC:

402          Hydrologic Process Parameterization of Stochastic Geoelectrical Imaging of Solute

403          Plumes: *SEG Technical Program Extended Abstract*, October 14 - 18, Anaheim, CA,

404          Proceedings, 5472–5476. 10.1190/segam2018-w12-01.1.

405 Oware, E. K., J. Irving, and T. Hermans (in-press), Basis-constrained Bayesian McMC

406          difference inversion for geoeletrical monitoring of hydrogeological processes,
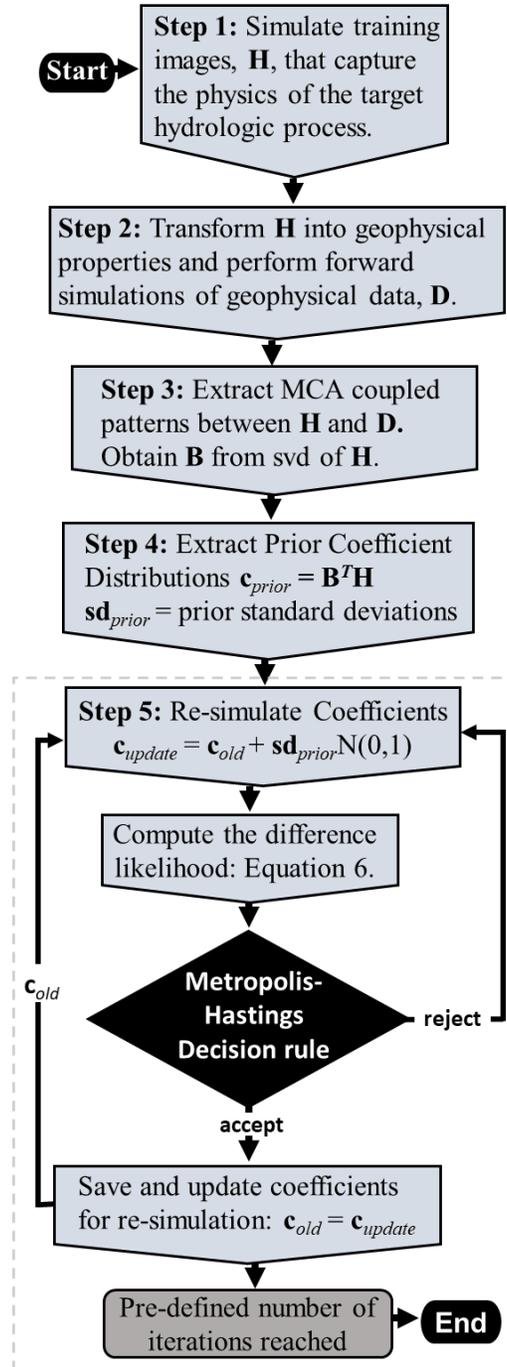
407          *Geophysics Letters*.

408 Satija, A., and J. Caers (2015), Direct forecasting of subsurface flow response from non-linear

409          dynamic data by linear least-squares in canonical functional principal component space

410          *Advances in Water Research*, 77, 69-81.

411 Scheidt C, P. Renard, and J. Caers (2015) Prediction-focused subsurface modeling: investigating

412          the need for accuracy in flow-based inverse modeling, *Math Geoscience*, 47, 73‑91.

413          http://dx.doi.org/10.1007/s11004-014-9521-6.

414 Tarantola, A. (2005), Inverse problem theory and methods for model parameter estimation,

415          *Society of Industrial and Applied Mathematics*.

416 Tippett, M. K., T. Delsole, S. J. Mason, and A. G. Barnstone (2008), Regression-Based Methods

417          for Finding Coupled, *J. Climate*, 21, 4384 - 4398.

418    von Storch, H. and F.W. Zwiers (1998), Statistical Analysis in Climate Research, *Cambridge*

419        *University Press, Cambridge*.

420    Zhou, H., J. J. Gómez-Hernández, and L. Liangping (2014), Inverse methods in hydrogeology:

421        Evolution and recent trends, *Advances in Water Resources*, 63, 22-37.
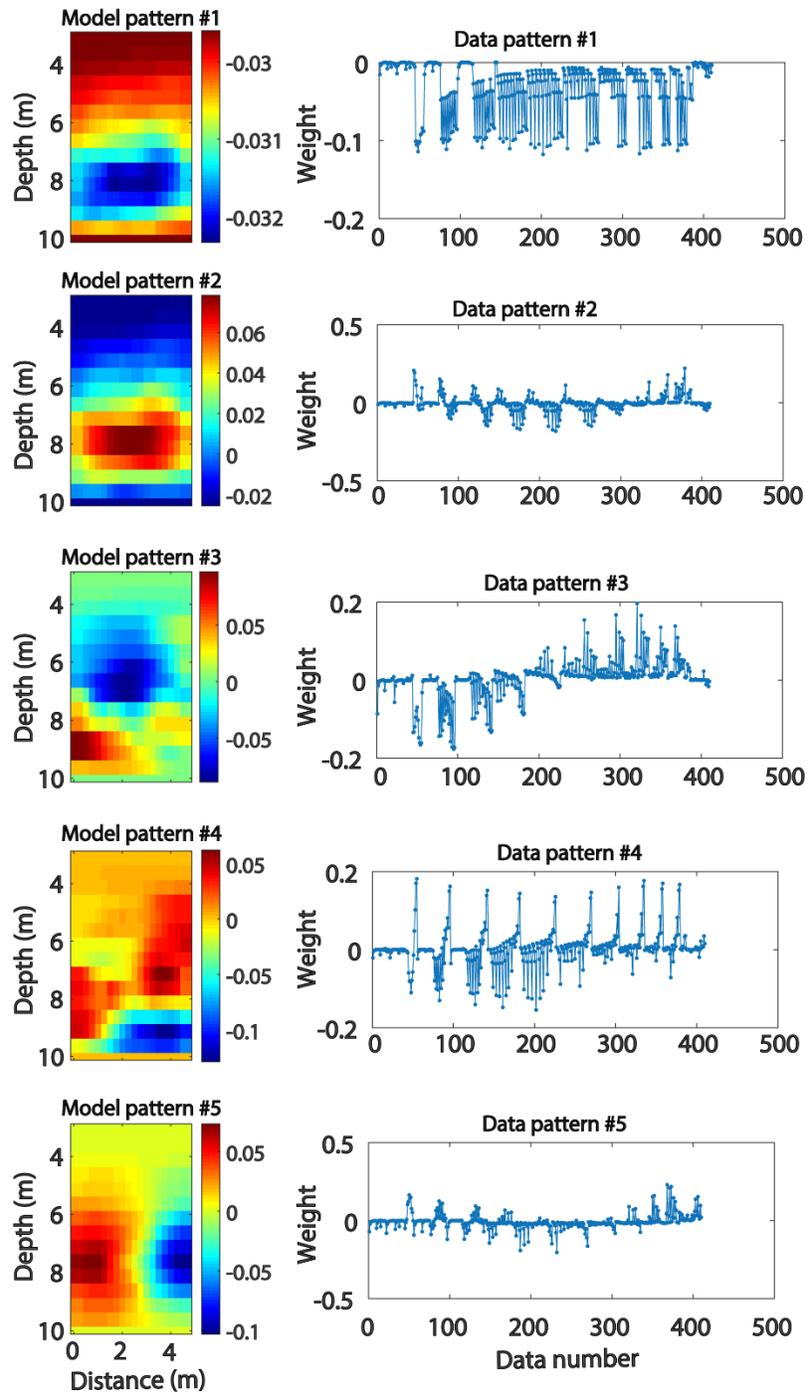
422

Figure 1: Flowchart for posterior sampling of the difference maximum covariance analysis prediction-focused approach (DMCA-PFA).
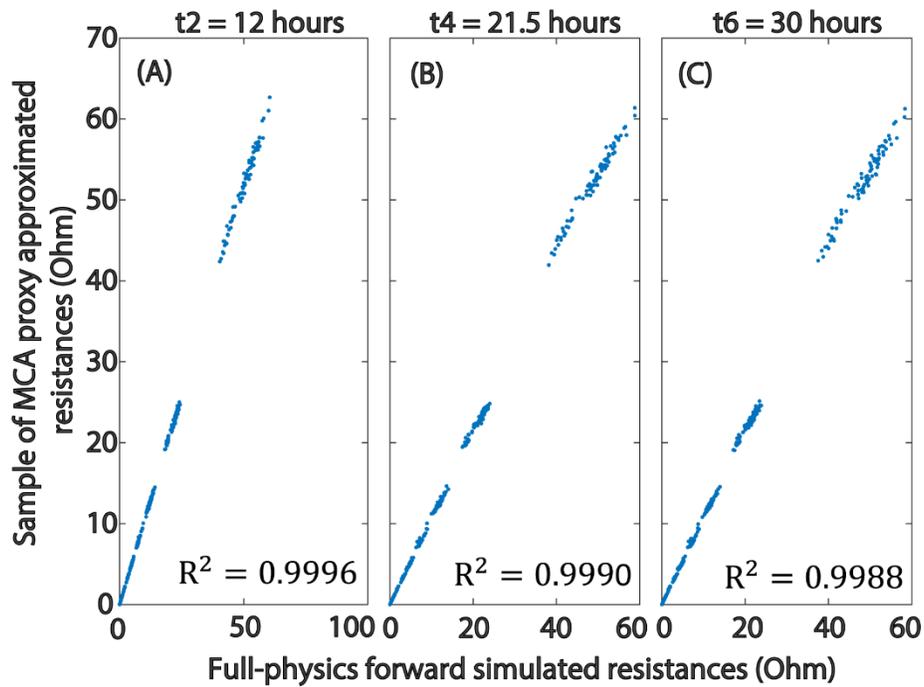
428

Figure 2: First 5 dominant maximum covariance analysis (MCA) coupled patterns between the
hydrological parameter (log(temperature); column 1) and geophysical data (resistivity; column
2) fields constructed from training images of 3000 mutual observations between the two fields.
The rows represent corresponding coupled patterns.

433
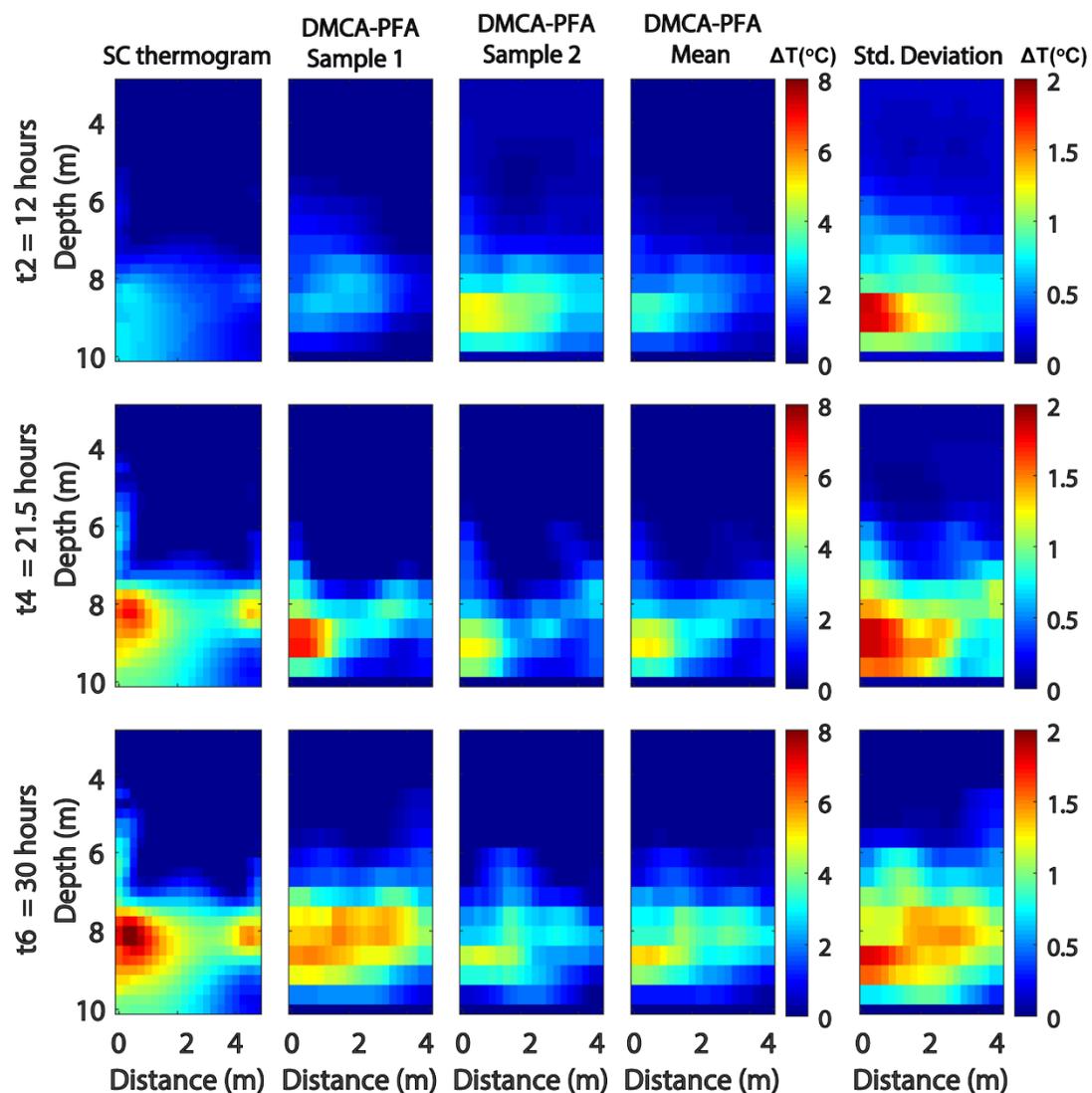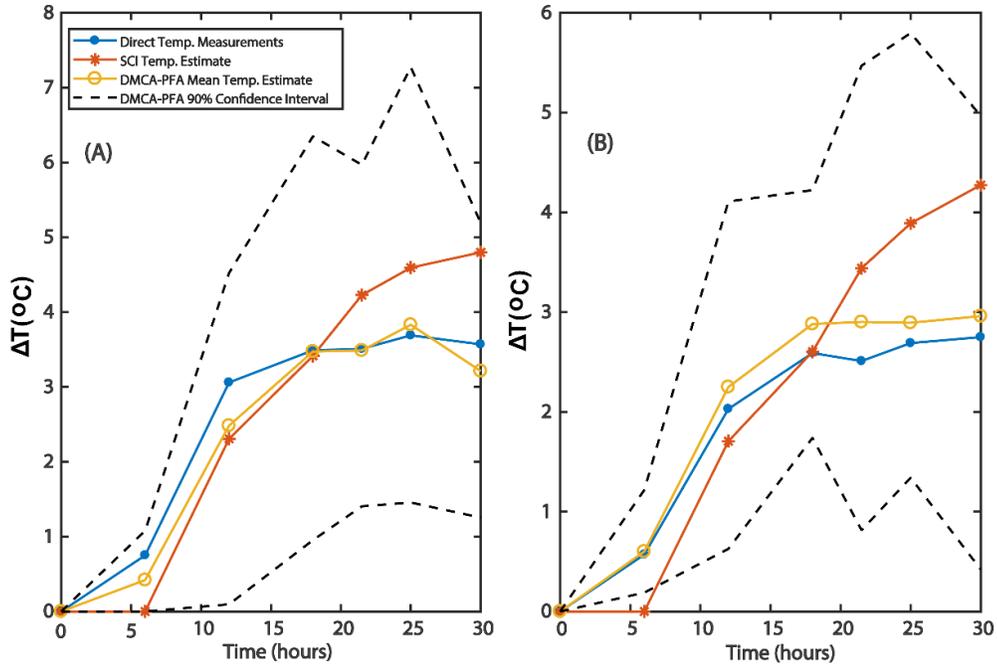
21

434

Figure 3: Scatter plots of high-fidelity resistivity forward simulated resistances vs sample of

435

maximum covariance analysis (MCA) approximated resistances obtained from smoothness-

436

constrained resistivity tomograms for time-step data at: (A) t2 (12 hours), (B) t4 (21.5 hours),

437

and (C) t6 (30 hours). The coefficients of determination ($R^2$) indicate almost a perfect one-to-one

438

MCA proxy-approximation of the high-fidelity forward simulated resistances.

439

440

441

442



Figure 4: Difference thermograms recovered directly from the ERT measurements at three different time steps: (row 1) 12h, (row 2) 21.5h, and (row 3) 30h. Column 1 shows thermograms from the classical smoothness-constrained (SC) inversion, columns 2 , 3, 4, and 5 show, respectively, two realizations, posterior mean and standard deviations estimated from the difference maximum covariance analysis prediction-focused approach (DMCA-PFA). Piezometers pz14 and pz15 are, respectively, located at (1.125 m, 9 m) and (2.25 m, 8.5 m).

Figure 5: Validation of estimated temperature break-through curves at two validation locations: (A) pz14 and (B) pz15. (Blue lines) direct temperature measurements, and estimated temperature break-through curves from the: (orange lines) classical smoothness-constrained inversion (SCI), (yellow lines) posterior mean of the difference maximum covariance analysis prediction-focused approach (DMCA-PFA) estimates. The two black dashed lines define the 90% confidence interval of the DMCA-PFA predictions.