1 **Predicting Future Well Performance for Environmental Remediation Design using Deep**

2 **Learning**

3 **Authors:** Xuehang Song[a,*], Huiying Ren[a], Zhangshuan Hou[a], Xinming Lin[a], Marinko

4 Karanovic[b], Matt Tonkin[b], Vicky L. Freedman[a,c], Inci Demirkanli [a,*], and Rob Mackley[a]

5 **Affiliations:**

6 [a] Pacific Northwest National Laboratory, 902 Battelle Boulevard, Richland, WA 99354

7 [b] S.S. Papadopulos and Associates Inc., 505 N. Pine St., Williamsfield, IL, 61489

8 [c] Sealaska, 1200 6th Ave, Suite 800, Seattle, WA 98101

9 *Corresponding Authors: xuehang.song@pnnl.gov and inci.demirkanli@pnnl.gov

10

11 **Abstract**

12 In this study, we developed a deep learning (DL) framework with a multi-channel three-

13 dimensional convolutional neural network (MC3D-CNN) to predict well performance and

14 thereby assist future environmental remediation design. Such prediction of extraction well

15 performance at designated locations is critical for configuring pump-and-treat (P&T) well

16 network design and operation, setting reasonable target closure dates for overall remedying, and

17 estimating remedy costs. The framework is developed with operational and monitoring data

18 routinely collected during P&T remedy operations, including well extraction and injection rates

19 as well as in situ contaminant concentrations. Traditionally, the collected data were rarely used

20 for purposes other than assessing **past well performance** and the accuracy of the conceptual site

21 model. However, recent advances in data-driven computational approaches enable better use of

22  the large datasets to inform **future well performance**, enhance site characterization, and

23  improve remediation planning. In this study, we established a DL framework to integrate

24  transient three-dimensional contaminant plumes and multiple aquifer properties (e.g., hydraulic

25  conductivity and hydrostratigraphic maps) to identify characteristic patterns controlling and

26  representing extraction well mass recovery, aiming at providing future mass recovery estimates

27  for existing wells and candidate wells at any proposed locations. We evaluated our framework by

28  using a realistic synthetic dataset generated from a well-calibrated flow and transport model used

29  in the 200 West Area of the U.S. Department of Energy's Hanford Site in southeastern

30  Washington state. The multi-channel feature in our framework allows integration of various

31  types and temporal densities of training datasets for DL model development. Overall, we found

32  that the trained DL model achieved an accuracy of over 90% in ranking extraction well

33  performance in validation datasets, and over 80% in predicting high-performance-ranking well

34  locations. This data-informed approach provides a flexible tool to support adaptive site

35  management, streamline decision-making, and potentially reduce remediation time and costs.

36  Our DL framework can be used as a filtering tool to improve the current P&T network

37  optimization design by reducing the number of candidate well locations.

38

43

## 1. Introduction

In this study, we detail a new framework for improving prediction of future well performance by using a deep learning (DL) model to discover the relationship between historical pump-and-treat (P&T) records and related contamination distribution, hydrostratigraphic units, and hydraulic conductivity maps. P&T is one common approach used to hydraulically contain and remediate groundwater contaminant plumes at many waste sites (McKinney and Lin, 1996; National Research, 2013; Truex et al., 2017). Statistics show that P&T is employed in approximately 40% of contaminated groundwater sites (EPA, 2002). In a typical P&T system, the contaminated groundwater is extracted to the ground surface by pumping and then treated using a filtering or stripping system. The treated water is then injected into the aquifer for groundwater recharge and/or to hydraulically contain the remaining contaminant plumes (EPA, 2005). The efficiency of P&T systems heavily depends on the performance of extraction wells, i.e., how much contaminant mass can be pumped from the aquifer within a reasonable time frame. Knowing extraction well performance at designated locations is critical for planning and modifying P&T systems, setting reasonable target closure dates for remedy operation, and estimating remedy costs (Zheng and Wang, 2002). While large amounts of historical well performance data are often routinely recorded under regulatory requirements, such information is mainly used for monitoring purposes only.

Over the past four decades, a variety of methods have been developed to simulate/predict the performance of extraction wells for remedy design, including the classical "batch flushing" analytical equation (Haley et al., 1989; National Research, 1994), analytic element method (Gaur et al., 2011; Majumder and Eldho, 2016; Matott et al., 2006), boundary element method (Kontos

3

67  and Katsifarakis, 2017), semi-analytical solutions (Ameli and Craig, 2018; Cardiff et al., 2010),

68  and more complex numerical fate and transport (F&T) models (such as MODFLOW, eSTOMP,

69  PFLOTRAN, ITOUGH2, and many other simulators) (Finsterle, 2006; Finsterle and Zhang,

70  2011; Hammond and Lichtner, 2010; Neville and Tonkin, 2004; White and Oostrom, 2003).

71  Among these approaches, the numerical F&T models usually provide more comprehensive

72  representations of complex site features, such as three-dimensional nonuniform distribution of

73  the contaminant plume, heterogeneous aquifer hydrogeological properties, dynamic groundwater

74  gradient, and geochemical and/or biogeochemical reactions in the aquifer (Finsterle and Zhang,

75  2011; Huang and Mayer, 1997; Minsker et al., 2004). However, calibrating a numerical model is

76  not trivial and can be time-consuming and site-specific. In many cases, site measurements are

77  inadequate to fully constrain model parameters, which results in non-unique solutions and high

78  predictive uncertainties of numerical models (Carrera, 1993; Singh and Minsker, 2008; Wu and

79  Zeng, 2013). Reactive transport modeling is especially hard because of its large problem

80  dimension (e.g., number of reactants and possible pathways) and high computational costs

81  (Mayer et al., 2001; Steefel et al., 2015; Tsang et al., 2015). Because of these difficulties,

82  simplified well models are still widely used for remedy planning (EPA, 2005). In addition, one

83  commonly used phased strategy in remedy design is to move from simple calculations to

84  analytical models and finally to a more detailed P&T system design (McMahon et al., 2001). In

85  spite of the great value and wide application of simplified well models in remedy design, it is

86  widely acknowledged that the existing simplified models lack representation of some important

87  site features, especially the nonuniform distribution of contaminant plumes and heterogeneous

88  aquifer hydrogeological properties, which can lead to overestimation of extraction well

89  performance and then underestimation of cleanup time and cost (Brusseau, 1996; Hadley and

4

90   Newell, 2012; National Research, 1994). These over-simplifications can be alleviated by

91   adjusting empirical parameters and adding additional components to the analytical solutions to

92   mimic more complex site behaviors (Sváb et al., 2008).

93

94   Recently, machine learning (ML) techniques have drawn increasing attention in water and

95   geoscience research (Shen, 2018; Tahmasebi et al., 2020). One popular application of ML

96   methods in geoscience is to construct surrogate models to substitute for computationally

97   expensive F&T models (Razavi et al., 2012; Yadav et al., 2018). In such cases, hundreds or a

98   few thousand realizations of physical F&T were created by perturbing model parameters, well

99   configurations, and/or initial conditions, and then the simulated modeling results were used as

100  training datasets to feed into regression ML models. The trained ML models have several

101  advantages, including high accuracy in reproducing the F&T model training dataset, low

102  computational requirements relative to F&T models, and the capability to generate new

103  predictions almost instantaneously. Mainly due to their high efficiency, the ML-based surrogate

104  models became popular in computationally demanding F&T optimization applications. A variety

105  of ML methods, such as artificial neural networks (Gaur et al., 2013; Rogers and Dowla, 1994;

106  Yan and Minsker, 2006), extreme learning machine (Majumder and Lu, 2021), and deep neural

107  network models (Chen et al., 2022; Yu et al., 2020), have been used to train surrogate models

108  and coupled with global optimization models (e.g., evolution algorithms) for P&T network

109  design. To alleviate the potential bias of surrogate models, it was also proposed to combine

110  multiple types of ML models and form more reliable ensemble surrogates (Yin and Tsai, 2020;

111  Zounemat-Kermani et al., 2021). One drawback of these surrogating approaches is that they lack

112  physical representation and may produce physically unrealistic results (e.g., violating mass

113    conservation law), and there is significant interest in enforcing physical laws to these "black box"

114    models through physics-informed ML methods, such as employing metamodels (Soriano et al.,

115    2021) or minimizing the residual of physical equations (Tartakovsky et al., 2020; Wang et al.,

116    2021), among others. In addition to being employed as emulators for F&T simulations, ML

117    methods are also used as novel inverse models in subsurface science by creating bidirectional

118    mappings between physical parameters (e.g., hydraulic conductivity and facies structures) and

119    model state variables (e.g., hydraulic head) (Mo et al., 2019; Sun, 2018; Wang et al., 2021). ML

120    models have also been used to improve monitoring and site characterization of active and closed

121    P&T sites by optimizing monitoring network design (Kontos et al., 2022; Meray et al., 2022),

122    improving plume source identification (Kontos et al., 2022), and filling data gaps (Ren et al.,

123    2022).

124

125    One recent intriguing topic of ML application in groundwater remediation is the use of pure

126    data-informed approaches. Instead of explicitly implementing physical constraints in ML models

127    or learning from physical F&T simulation results, these data-informed approaches seek to

128    directly link the model outcome of interest and its controlling factors under ML frameworks for

129    better groundwater contaminant estimation and remediation design. The selection of the

130    controlling factors is not arbitrary, but rather is based on the understanding of causal relations

131    of the aquifer system and physical laws. Thus, the ML algorithms are performed as a way of data

132    mining to extract and formulate hidden relations and patterns between the variables of concern

133    and their controlling factors. A predictive model of groundwater nitrate pollution was built using

134    random forest (R.F.) regression by examining nitrate concentrations with 24 related site

135    parameters, such as intrinsic hydrogeologic properties, driving forces, remotely sensed variables,

136    and physical-chemical variables (Rodriguez-Galiano et al. (2014)). McConnell and others trained

137    regression-based ML models (Prophet model and damped Holt's exponential smoothing model)

138    to predict future carbon tetrachloride ($CCl_4$) plumes using historical $CCl_4$ concentration samples

139    (McConnell et al., 2022). Their models can achieve satisfactory prediction of site closure time

140    without solving governing transport equations with sufficient spatial and temporal density of

141    data. Wu and others applied R.F. methods for classifying high health risk areas using various

142    groundwater chemistry measurements, and demonstrated that an R.F. model with four types of

143    the most important chemistry measurements can achieve a classification accuracy of 88.21% for

144    groundwater quality (Wu et al., 2020). The aforementioned applications focused on estimating

145    the extent of a current contaminant plume and/or predicting future plume migration; no studies

146    has been reported that predict the performance of the P&T system using data driven-approaches.

147

148    With the rapid advance in ML applications in contaminant migration modeling and remediation

149    design, one important dataset, historical P&T records, has been under-utilized. Under regulatory

150    requirements, large amounts of extraction and injection records, such as flow rate and

151    contaminant concentration, are routinely collected as a standard practice in most P&T

152    remediation sites. These data are mainly used to monitor contaminant rate/mass removal while

153    they directly reflect the response of the aquifer, and potentially can be mined to support better

154    decision-making for future remedies (Brusseau, 2013; Truex et al., 2017). On the contrary, the

155    historical well production data have already been proven to be very useful in the gas and oil

156    industry to predict future oil production rates and guide future well drilling; e.g., (Hirschmiller et

157    al., 2019); (Li et al., 2019).

158

159 Our DL approach allows integration of multi-type multi-resolution P&T monitoring and site

160 characterization data, addresses the limitations of the analytical and semi-analytical solutions in

161 representing heterogenous field characteristics, and avoids solving expensive governing transport

162 equations. An image-based DL model, convolutional neural networks (CNNs), was developed

163 and integrated to extract the hidden spatiotemporal correlations between physical control factors

164 and historical records of P&T well performance. The model is evaluated using a realistic

165 synthetic dataset generated from a calibrated F&T model for a site with historical contamination.

166 The effects of training dataset type and temporal data density are also interrogated to understand

167 potential improvements in model performance. Given the large amount of P&T records

168 generated at many waste sites for monitoring purposes, and the increasing automatic collection

169 and digitalization of these records, data-informed approaches such as our multi-channel CNN

170 create opportunities to improve our understanding of contaminant transport and site

171 management, streamline decision-making, and potentially reduce future remediation costs.

172

173 **2. Methods**

174 We developed a multi-channel 3D-CNN (MC3D-CNN) DL architecture to extract important

175 features from transient plume distributions and aquifer hydrogeological properties that control

176 extraction well performance. The trained DL model can be used as a prediction model that

177 provides favorable locations for future wells to maximize contaminant mass recovery, shorten

178 the operational time of the P&T system, and reduce total remediation cost. Section 2.1 introduces

179 the DL background and key configurations of CNN, and Section 2.2 describes the architecture of

180 the MC3D-CNN model (Figure 1) and its major hyperparameters. Implementation of the CNN

8

181    classification model for predicting well performance with physical model simulations and

182    physical properties is illustrated in Figure 2 and discussed in Section 2.3.

183

184    **2.1. Deep learning and CNN configurations**

185    DL is a sub-field of ML, which has been designed to reveal the hidden controlling mechanism in

186    high-dimensional and nonlinear complex systems. Typical DL approaches, such as fully

187    connected neural networks (FC-NNs), CNN, and long short-term memory, can automatically

188    find the most salient features to be learned. These approaches have demonstrated tremendous

189    success in a variety of applications, such as speech recognition, computer vision, and natural

190    language processing CNN has outperformed other DL methods in predictive capability in many

191    image-related applications, including medical imaging, material structure, object recognition,

192    and others (Rao and Liu, 2020).

193    CNN was first developed for visual imagery analysis and feature extraction (LeCun et al., 2015).

194    The Visual Geometry Group (VGG) block-wise model architecture is adopted to push the model

195    depth toward high accuracy (Simonyan and Zisserman, 2015). The VGG model architecture

196    includes a series of convolutional blocks containing multiple convolutional layers followed by

197    batch normalization, pooling, and dropout layers within each block, and then connected to flatten

198    and dense layers. VGG model architecture needs to be adjusted according to the characteristics

199    of the datasets. The number of convolutional filters is incremented layer by layer to make sure

200    that the increasingly richer features are properly extracted. Such layered organizations can learn

201    hierarchical representations. The neurons of adjacent layers are connected by assigning weights

202    and biases $\{W_i, b_i\}_{i=1}^{m}$, where $m$ is the number of layers in neural network $NN_m$. The initial layer

203    is the input layer constructed in image sets and the last layer is the output defined as the

204   classification labels. The predicted outcome is compared with the label and a measure is

205   calculated representing the performance of the CNN. The categorical cross-entropy class is

206   chosen for the multi-label classification problems. It computes the cross-entropy loss between

207   the labels and model predictions, and the calculation of the loss function requires that the last

208   dense layer be configured with the total number of classes; this enables softmax activation to

209   predict the probability for each class. In between are hidden layers transforming the feature space

210   of the input such that it matches the output. Max pooling is performed to reduce the data size

211   using spatial down-sampling, while preserving discriminant information. The normalization of

212   output from previous layers allows the neural network to learn the pattern more independently.

213   Dropout as a common regularization technique is also used to introduce stochasticity to make

214   model performance more robust and prevent overfitting.

215

### 2.2  Multi-channel 3D-CNN architecture

217   Three-dimensional (3D)-CNN is needed to take labeled 3D images for feature extraction, but it is

218   computationally and memory exhausting because of its much larger number of trainable

219   parameters compared to the regular two-dimensional (2D)-CNN variant. Recent advances in

220   computational hardware, especially general-purpose graphics processing units, have made 3D-

221   CNN computationally affordable (Zhao et al., 2019). In a typical 3D-CNN design, a 3D image

222   passes through a series of blocks of convolutional layers to extract feature maps, as shown by the

223   architecture illustrated in Figure 1. Multiple 3D image datasets, including different data types

224   and their temporal and spatial components, are fed into our multi-channel 3D-CNN. Under such

225   architecture, the ensemble of sub-CNNs per channel are trained simultaneously to learn the

226   spatiotemporal features ingested from various sources (e.g., plume distribution,

227    hydrostratigraphic unit map) and match the features to the predefined labels (e.g., well

228    performance ranking index).

229

230    Hyperparameter searching is needed to optimize the MC3D-CNN model configuration. A series

231    of configuration parameters were explored, including batch size, kernel size, number of layers,

232    number of neurons in each convolutional block, and dropout rate. The optimal configuration was

233    then chosen by comparing the performance metrics of various hyperparameter combinations on

234    training and validation datasets. The final MC3D-CNN model was then evaluated with an
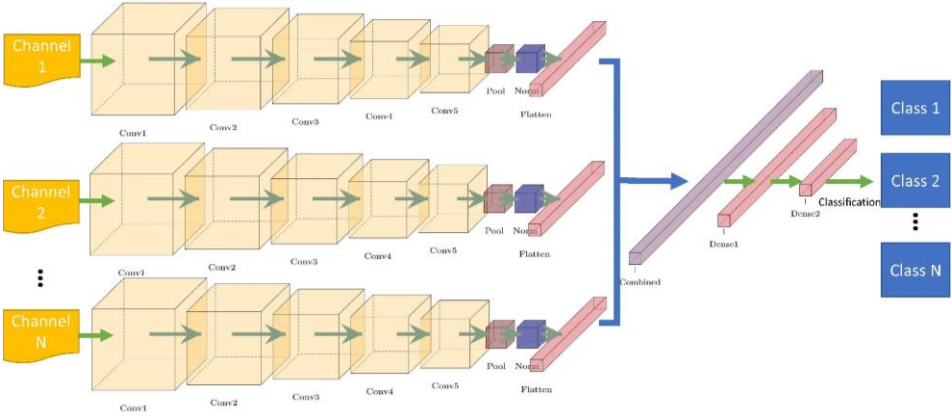
235    independent testing dataset.

236

237

238



239                    Figure 1. Architecture of the MC3D-CNN classification model.

240

241    **2.3 DL framework for well performance classification**

11

242 Figure 2 illustrates the overall workflow for training MC3D-CNN and predicting future

243 performance ranking. Three physical attributes: plume concentration, hydrostratigraphic unit,

244 and hydraulic conductivity, were used to compose the 3D training images. In this study, a

245 synthetic modeling dataset was used to generate the 3D training images for each extraction well,

246 as detailed in section 3.2.2. The accumulative mass recovery was categorized into three levels of

247 performance, which serve as the image labels for CNN supervised learning. The raw pixel data

248 from the 3D training images were fed into the 3D-CNN model, which can integrate various

249 sources of data representing different aspects of system behaviors and dynamics related to well
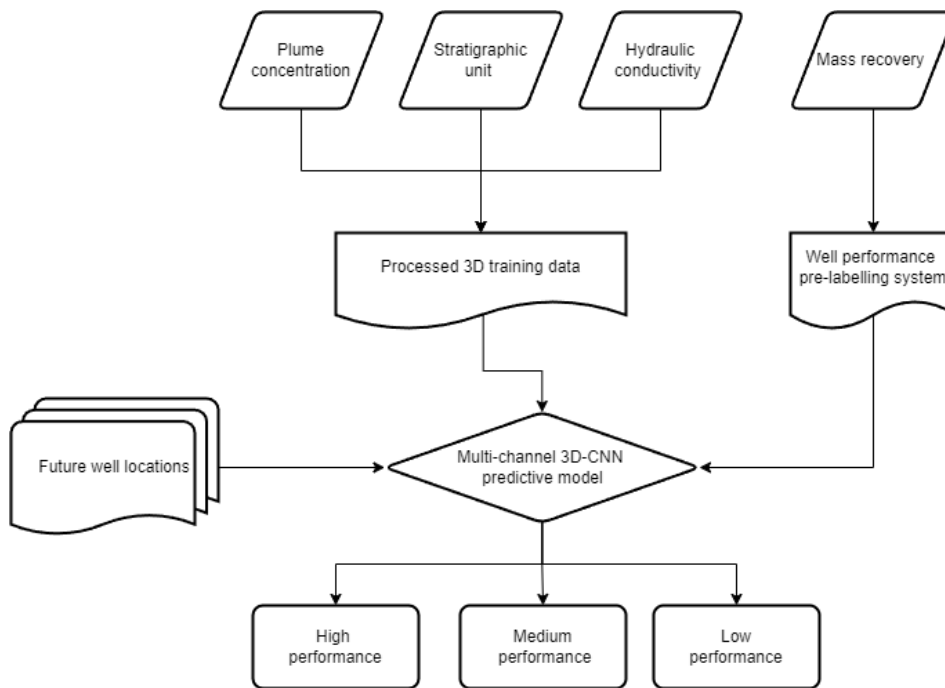
250 performance.

251



252

253         Figure 2. DL-based workflow for ranking well performance.

254
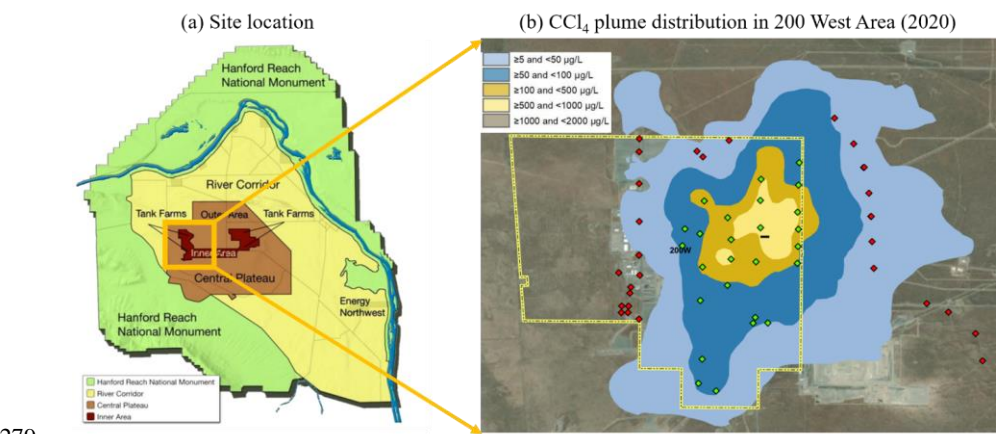
## 3. Case Design and Model Evaluation Criteria

256    A suitable example for demonstrating the MC3D-CNN model is the remediation of the Central

257    Plateau at the U.S. Department of Energy (DOE) Hanford Site (Section 3.1). Here, we extracted

258    the simulated $CCl_4$ plumes from an F&T model covering the Central Plateau (Section 3.2) to

259    compose a synthetic dataset for model testing. The advantage of using synthetic data from the

260    F&T model is that it can provide a known reference for evaluating the accuracy of the DL model.

261    Section 3.3 presents the details of the DL model setup and model parameters for the synthetic

262    case. The criteria used to measure the performance of the DL model are described in Section 3.4.

263    Section 3.5 illustrates how to use the trained DL model to predict well performance ranking.

264

## 3.1. Site description

266    The DOE Hanford Site, located in southeastern Washington State, holds radioactive waste from

267    the disposal of nuclear fuel fabrication wastes from 1943 to 1975 (Figure 3). The Central Plateau

268    is an informal geographic designation given to the broad central portion of the Hanford Site that

269    encompasses the 200 West and 200 East Areas. The Central Platea area is one of the most

270    complex environmental remediation sites in the world, with shallow sources (e.g., waste tanks),

271    persisting and recalcitrant deep vadose zone residual sources, large-scale groundwater plumes

272    [e.g., carbon tetrachloride ($CCl_4$), technetium-99 (Tc-99), iodine-129 (I-129) and nitrate ($NO_3$)],

273    and subsurface heterogeneity (Demirkanli and Freedman, 2021). Groundwater with several

274    contaminants of concern (COCs) has been treated by the 200 West P&T Facility in the Central

13

275     Plateau since 2012. The 200 West P&T Facility is designed to capture and treat contaminated

276     groundwater to reduce the mass of selected COCs, such as $CCl_4$, $NO_3$, Tc-99, and I-129, by at

277     least 95% within 25 years from the startup (Demirkanli et al., 2018). The locations of extraction

278     wells are shown as green diamonds in Figure 3b.



279

280     Figure 3. (a) Site location (source:

281     https://www.hanford.gov/files.cfm/Attachment_5_Approach_CP_Cleanup_handout.pdf ); (b)

282     $CCl_4$ plume distribution in 200 West Area (downloaded from

283     https://www.hanford.gov/page.cfm/PHOENIX). The yellow polygon shows the boundary of the

284     200 West Area. The green diamonds are existing extraction wells and the red diamonds are

285     existing injection wells.

286

287     **3.2. Synthetic dataset generation**

288     The Plateau-to-River groundwater model (P2R model) (Budge and Nichols, 2020) has been

289     developed for the Central Plateau and extends eastward to the Columbia River (Figure S1). The

14

290  P2R model was calibrated using hundreds of monitoring wells. The model primarily provides the

291  computational basis for simulating the F&T of contaminants in groundwater within the near- and

292  far-field portion of the affected aquifer in the Central Plateau, and is currently used to support

293  ongoing remedial activities on the Central Plateau. More details of the P2R model can be found

294  elsewhere (Budge and Nichols, 2020).

295

296  **3.2.1. Well performance ranking**

297  While the P2R model simulated groundwater F&T of multiple contaminants in the Central

298  Plateau, this study only focused on the $CCl_4$ removal data in the 200 West Area for

299  demonstration purposes. The performance of each well is ranked as high, medium, or low

300  according to the well's $CCl_4$ recovery. These well rankings were used as labels in the DL

301  classification model. It is noticed that the ranking of each well varies over time, where a typical

302  well might be ranked as a high-performance well in its early years of operation, then its

303  performance will decrease over the years with the removal of surrounding $CCl_4$. Therefore,

304  performance of each well is distinguished and labeled at multiple time segments. Figure 4 shows

305  the simulated annual $CCl_4$ mass recovery of 28 existing extraction wells from the P2R model.
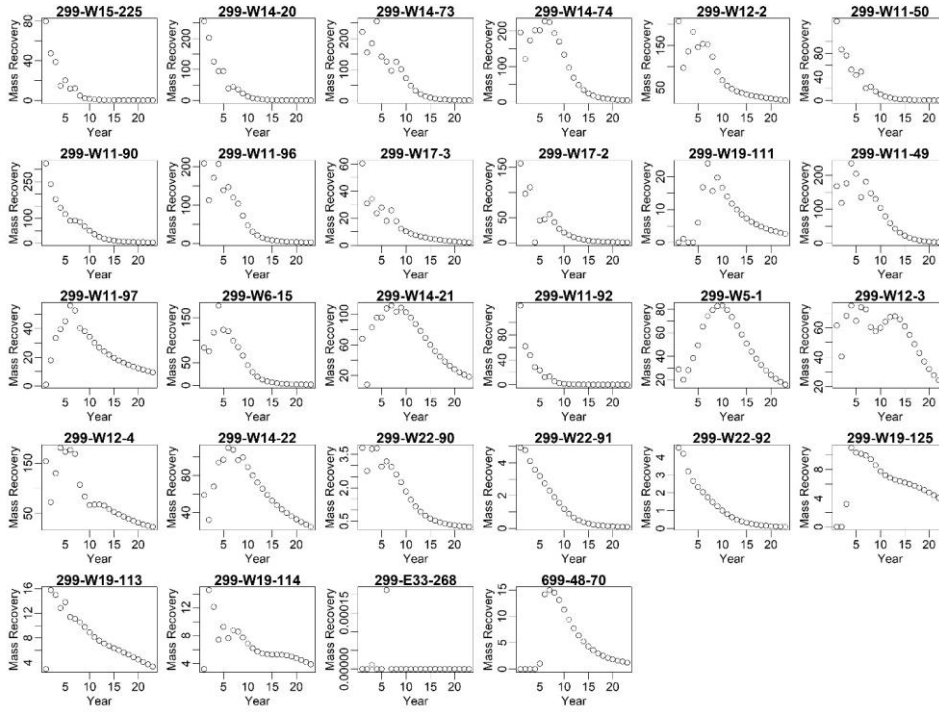
15

306

Figure 4. Yearly mass recovery of 28 wells extracted from P2R model simulation.

307

308

309    Each of the 28 extraction wells was then labeled as high, medium, or low performance according

310    to its 5-year cumulative mass recovery (Eq 1), as:

$$R_{Low}: M_{i,j} < C_1$$

(1)

$$R_{Medium}: C_1 \leq M_{i,j} \leq C_2$$

$$R_{High}: C_2 \leq M_{i,j}$$

16

311

where $C_1$[Kg] and $C_2$[Kg] are predefined threshold values. $R_{Low}$ , $R_{Medium}$, and $R_{High}$ denote the

performance indicator values corresponding to low, medium, and high performance,

respectively. $M_{i,j}$ is a moving sum calculated from the annual CCl₄ mass recovery:

$$M_{i,j} = \sum_{k=j+1}^{j+n} m_{i,k} \text{ , for } j = N - n \qquad (2)$$

315

where $m_{i,k}$[Kg] is the CCl₄ mass recovery of well $i$ in year $k$ ; $M_{i,j}$ is the CCl₄ mass recovery of

well $i$ for the following $n$ year starting in year $j$; $n$[year] is the time window for moving sum,

which is 5 years in this study; and $N$[year] is total number of years. For example, the "future"

performance rank of a well in 2012 is evaluated using its total mass recovery between 2012 and

2017. We chose a 5-year moving time window because more than half of the wells reached their

peak performance around year 5 (Figure 4). The values of C1 and C2 are 30[Kg] and 200[Kg],

respectively, where the selections are based on the quantile and distribution of the cumulative

mass recovery and the pre-designed equal distributed number of members in each class (Figure

S2). The purpose is to make a balanced classification system to improve the predictive

performance and avoid model bias. Figure 5 shows the calculated ranking for each well in

different years. It is not surprising that most of the 28 wells move from medium/high

performance to low/medium performance rankings over decades of remediation. However, there

are still a few wells that move from low/medium to medium/high rankings (e.g., 299-W11-49

and 200-W11-97). This is because the plume spreading increases the CCl₄ concentration in those

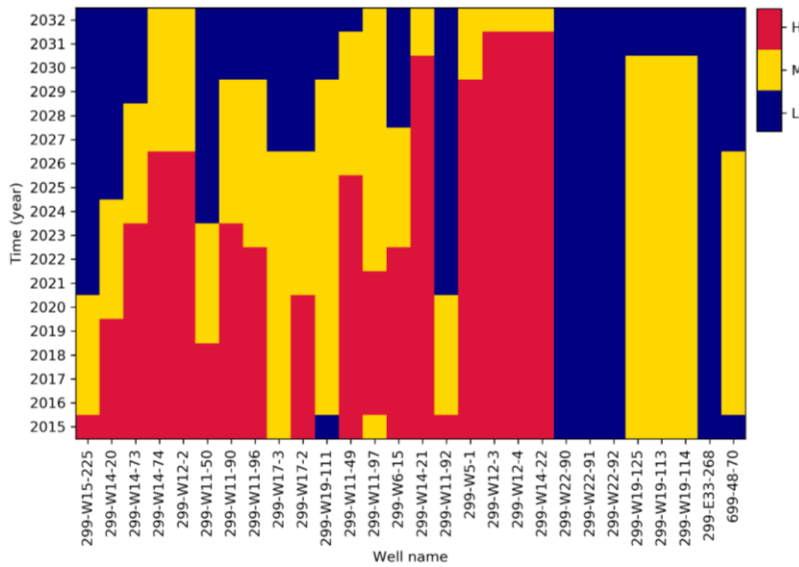wells and thus improves their CCl₄ extraction rates.

331

Figure 5. Well performance rankings for DL model training (year 2015~2028) and testing (year

2029 ~2032).

### 3.2.2. 3D image sets under different scenarios

Around each well, three types of model parameters and variables – plume concentration,

hydraulic conductivity, and hydrostratigraphic unit – were extracted from the P2R model

configuration and output files. The horizontal and vertical lengths of the 3D images are 1000 m

and 80 m, respectively, as selected by model sensitivity test. Three scenarios have been designed

to test different combinations of these datasets on DL model performance:

- Scenario 1 (S1) uses the current snapshot of plume concentration only.

- Scenario 2 (S2) adds two earlier time steps to S1 as the model inputs. Such a setup

    considers the temporal impact for model prediction.

18

343    • Scenario 3 (S3) adds hydraulic conductivity and hydrostratigraphic units to S2. Both

344        hydraulic conductivity and hydrostratigraphic units are static over time.

345    Examples of different data sources are illustrated in Figure 6 at three wells representing different

346    performance indicators. In the year 2021, the performance wells 299-W14-20, 299-W14-73, and

347    299-W-15-225 are ranked as medium, high, and low, respectively. In the year 2026, the

348    performance rankings of the three wells change to low, medium, and low, respectively. In the

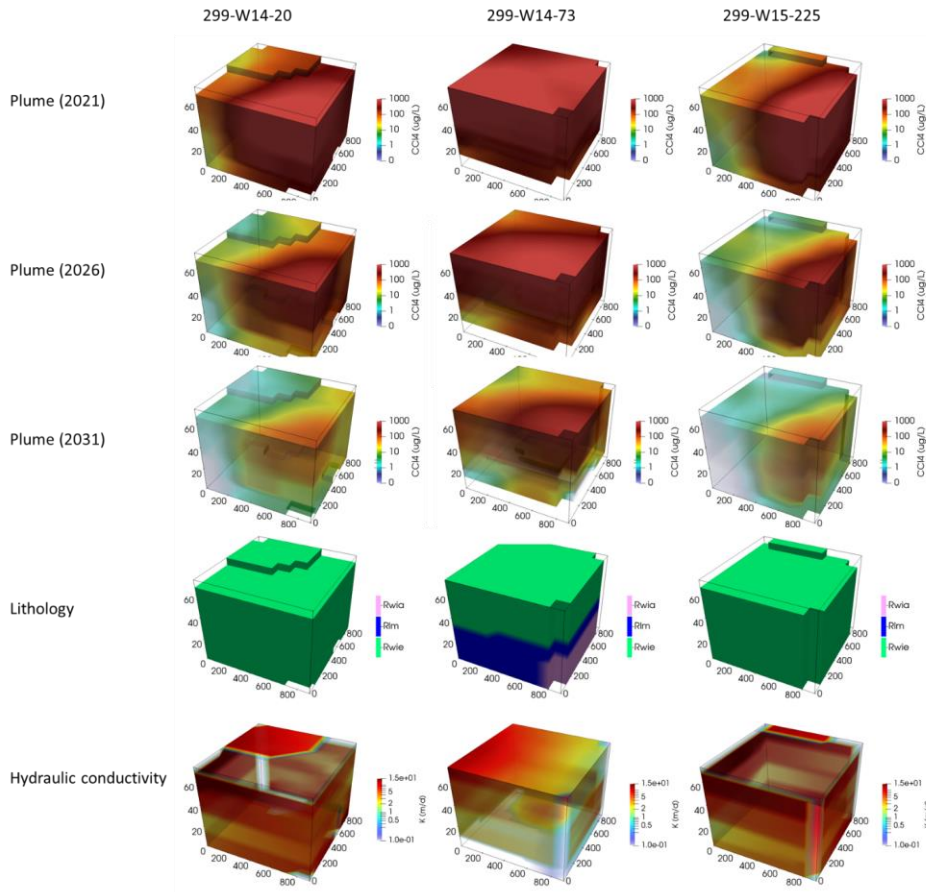349    year 2031, all three wells are ranked as low performance.

350

Figure 6 . Examples of 3D training image sets from multiple data sources at three different wells.

Each well is located at the center of the training images.

### 3.3. DL model setup

Based on the performance ranking rule defined in Section 3.1, a total of 486 pairs of label

(performance ranking) and training image datasets were created for the 28 wells at 18 timesteps.

Note that not all wells were operated throughout the entire simulation time frame. We first used

20

357   80% of the dataset for model training and validation (year 2015~2029) and the remaining 20% of

358   the dataset for model testing (year 2029~2032). Each scenario described in Section 3.2 shares the

359   same 3D-CNN model architecture and grid searching in the hyperparameter tuning process to

360   make a fair comparison across different scenarios.

361

362   **3.4. Model accuracy measures**

363   The accuracy of the trained ML model was evaluated in the testing datasets, which are the last

364   20% of the synthetic data (year 2029~2032). The direct model predictions were plotted in

365   corresponding years to examine misclassified testing points spatiotemporally. The commonly

366   used confusion matrix was applied to quantify the MC3D-CNN model performance using the

367   testing set. The confusion matrix visualizes the accuracy of a classifier by comparing the actual

368   and predicted classes (e.g., well performance rankings in this study). For a binary classification

369   problem, the four types of events in its confusion matrix and their meanings are as follows:

370       (1) True Positive (T.P.), correctly predicted true positives.

371       (2) False Negative (F.N.), true positives predicted as negative values.

372       (3) False Positive (F.P.), true negatives values predicted as positives

373       (4) True Negatives (T.N.), correctly predicted true negatives

374   For multi-class classification, T.P., F.N., F.P., and T.N. need to be determined for each class

375   separately by lumping all other classes into one class. Then, the confusion matrix was used to

376   assess the accuracy, precision, sensitivity, and specificity of each class to diagnose the model

377   performance. Accuracy represents the percentage of correctly labeled events within the whole

378   pool of events. Precision uses the portion between T.P. and (T.P.+F.P.) to assess how good the

379     model is at assigning positive events to the positive classes. Sensitivity is the fraction between

380     T.P. and (T.P.+F.N.), which measures how proper the model is for detecting events in the

381     positive class. Specificity is the ratio of T.N. and (T.N.+F.P.), which evaluates how exact the

382     assignment to the negative class is.

383

384     **3.5. Future well performance prediction**

385     The trained DL model was used to predict the future well performance ranking in year 2022 for

386     the entire model domain to demonstrate the usage of this model. The unconfined aquifer in the

387     200 West Area was first discretized to a structured grid with a spatial resolution of

388     $100 \times 100 \times 5\ m$. Imaginary wells were placed in each grid node with a predefined screen

389     length of 48 m, which is the medium value of the existing wells. Same as the training dataset, a

390     series of $10 \times 10 \times 13$ image sets were extracted from $CCl_4$ plume, hydraulic conductivity, and

391     hydrostratigraphic unit datasets as inputs for model prediction. The future well performance

392     ranking was then generated using the trained MC3D-CNN model for each of the gridded

393     locations, which were then used to create a map of well performance ranking.

394

395     **4. Results**

396     We trained DL models for the three designed model configurations, and in Section 4.1 we

397     evaluate their accuracy. The trained models achieved over 90% accuracy on the training and

398     validation datasets, and provided satisfactory results on the testing set. The trained DL models

399     were then used to provide a site-wide ranking map to illustrate the usage of this method (Section

400     4.2).

401

**4.1. Model evaluation under different scenarios**

The CNN model was trained with various configurations for well performance evaluation and

the optimal model configuration was applied to the third independent testing dataset. The model

training history and class statistics were calculated and are illustrated in Figure 7. Figure 7 (a-c)

represents the influence of the optimized model configurations (e.g., the model accuracy vs.

epochs curves) with three settings(scenarios) of predictors: single plume (S1), multi-step plume

(S2), and multi-step plume with field properties (S3). Both the training and validation accuracy

increased with epochs and converged at or above 90% for all three scenarios. With more data

channels added to the training pool, the model accuracy increased. Although the S3 model has

the highest training and validation accuracy, the overall averaged accuracy for all three models

was satisfactory without noticeable over-fitting. It is not surprising that S2 and S3 yielded similar

learning results because (1) plume distribution is the most important control factor on extraction

well performance and (2) the impacts of hydraulic conductivity and hydrostratigraphic units have

been implicitly represented in the multi-year plume variations.

416

Figure 7 (d-f) is a multi-class confusion matrix for the testing set; specifically, the diagonal line

stands for the matched cases between predictions and targets for each class, the upper corners are

the overestimated cases, and the lower corners are the underestimated cases. In general, the

results show that the testing accuracy across the three scenarios is lower than training and

validation accuracy. The reason is that the model was trained using 12 years of historical data to

predict 4 years of future behaviors. This indicates that the aquifer conditions were nonstationary

due to the continuous decrease of $CCl_4$ inventory. For the single-plume scenario (S1), the overall
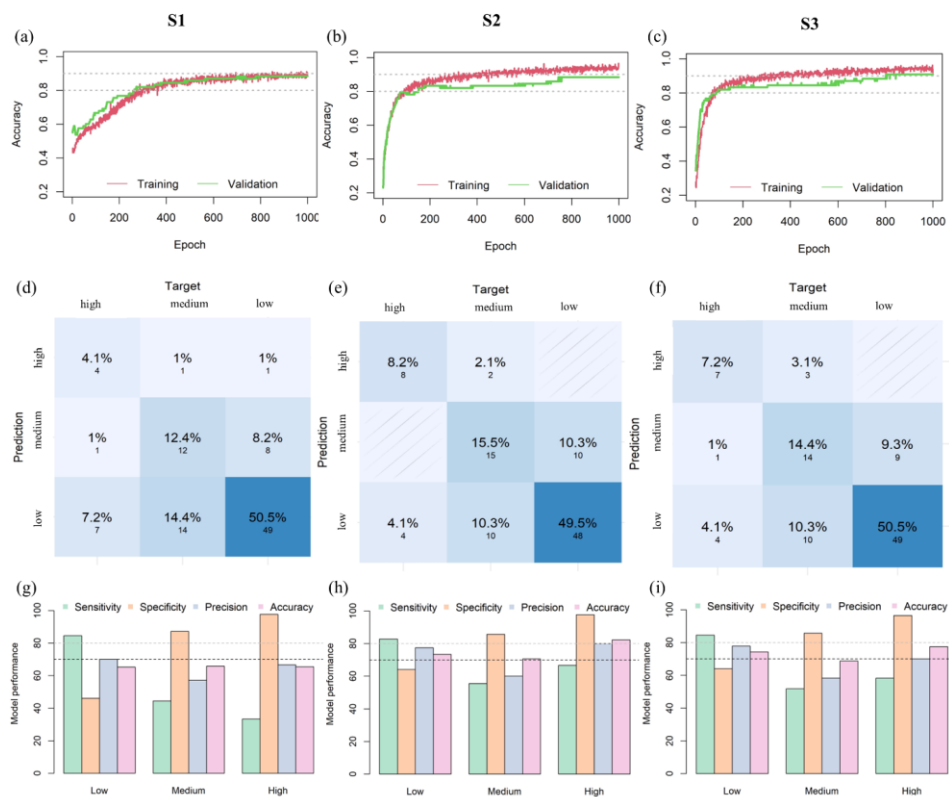
23

424  model accuracy for the testing dataset is 67%. The model has about 23% underestimation on well

425  performance, where 7 high-performance-ranking cases and 14 medium-performance-ranking

426  cases were classified as low-performance-ranking. The model has good performance in

427  controlling the overestimation, with only 8 low-performance-ranking cases classified as medium.

428  After adding multi-time channel data (S2 and S3), the averaged model accuracy increases to 73%

429  and the improvement is observed for both medium and high classes. The underestimated rate was

430  reduced to 14%, which is 8.6% lower than S1. The overestimated rate was 2% higher compared

431  to the single-plume scenario. The DL models trained by S2 and S3 have comparable

432  performance.

433

434  Figure 7 (g-h) shows class statistics using four metrics calculated from the confusion matrix. For

435  the single-plume scenario (S1), the sensitivity for the high and medium classes is low because

436  the number of T.P. for both classes is small, which means the model needs to be improved to

437  better detect high- and medium-performance cases. The specificity of the low class is 46%,

438  which means the model cannot assign low class exactly; in this case, medium- and high-

439  performance-ranking cases are likely to be predicted as low performance. All four performance

440  statistics improve after adding multi-step temporal information (S2 model). In the high-

441  performance-ranking class of the S2 model, accuracy and precision reached 82% and 80%,

442  respectively. The high precision in predicting the high-performance-ranking class means that the

443  model is good at predicting high-performance-ranking cases. The S3 model demonstrates slightly

444  better sensitivity for the low-performance-ranking class, and its overall statistical performance is

445  similar to that of the S2 model. Based on the comparisons of the above statistical performance

24

446    metrics, S3 was selected and applied to the testing dataset and the field prediction described in

447    the following sections.



448

449    Figure 7. Model optimization and performance evaluation for different scenarios; each column

450    represents a pre-designed scenario: (a-c) model training history under the most suitable

451    configuration; (d-f) multi-class confusion matrix; (g-i) model performance metrics on the testing

452    set.

453

454    Figure 8 shows the spatial distributions of correct prediction, underestimation, and

455    overestimation cases by the S3 model. Each dot under the subplots represents a testing well

456    location for a particular year. In general, model predictions match the label references for most

457    wells. However, since the mass recovery data extracted from the P2R model tends to decrease

458    over time and the model predicts future well performance, the S3 model tends to underestimate

459    for later years. These results suggest that the deep learning model could be further improved by

460    including more representative spatiotemporal samples into the training set.
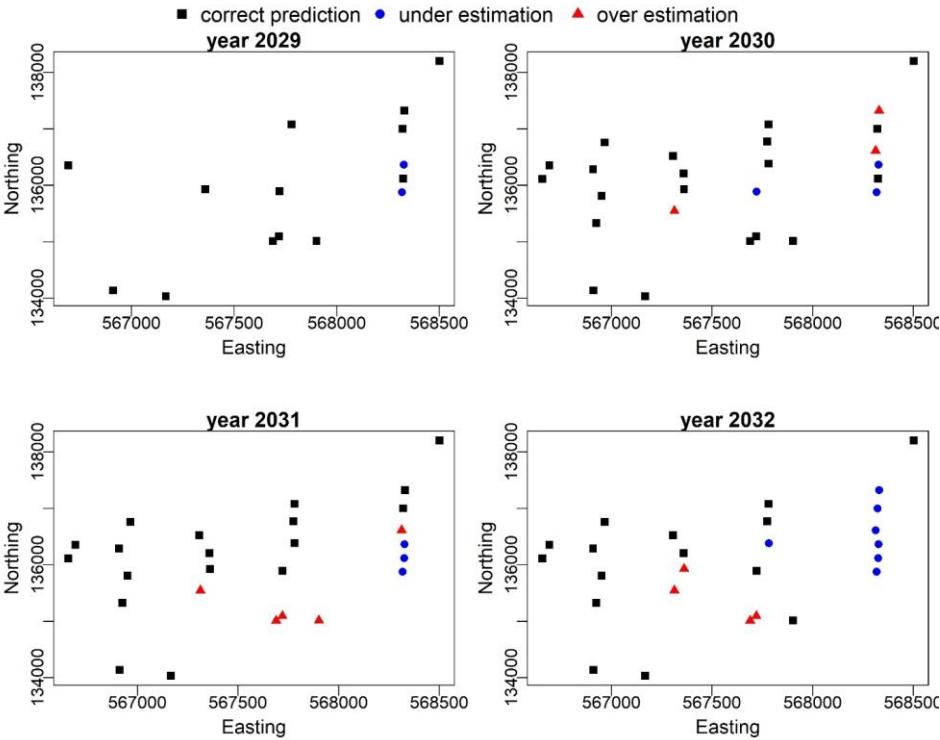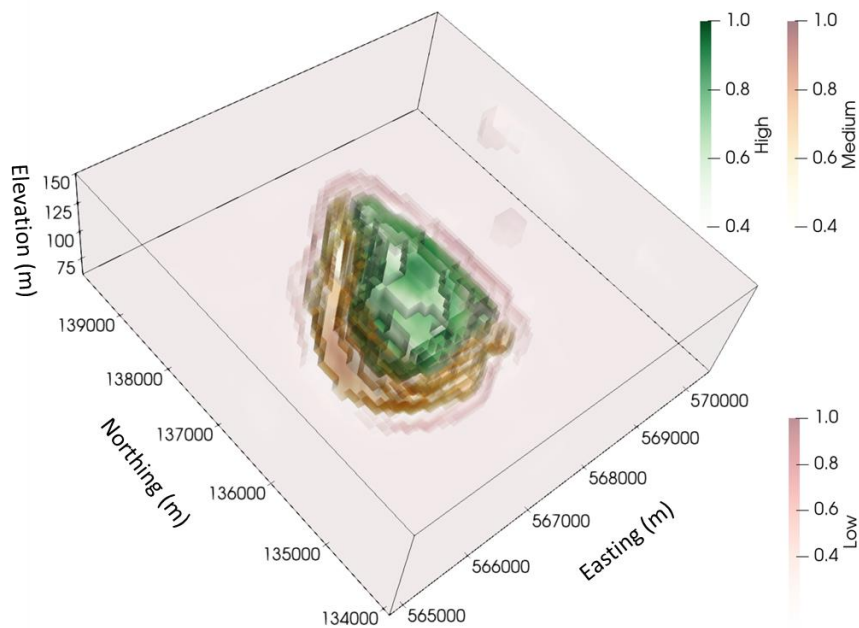
461



462

26

463  Figure 8. Spatial distribution of predictive well performance ranking obtained from the S3 model
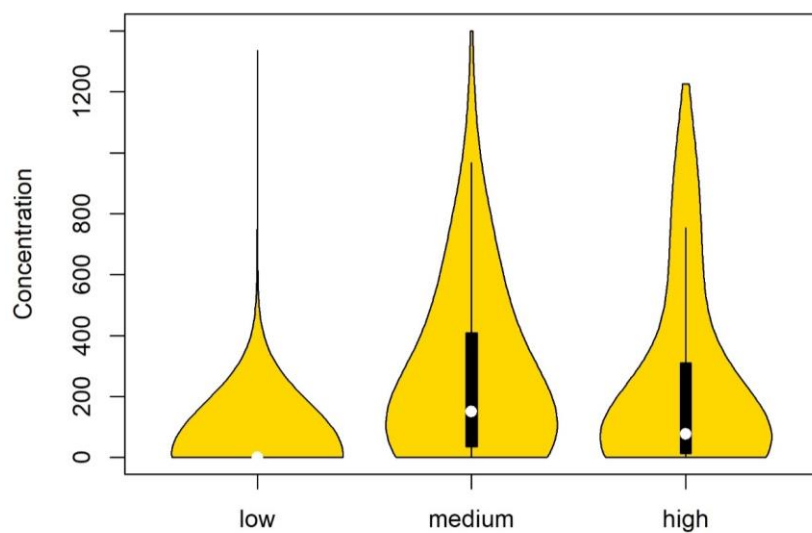
464  at the testing time period (year 2029~2032).

465

466  **4.2. Well performance ranking prediction**

467  The predicted well performance ranking map covering the entire field is obtained using the well-

468  trained model. Based on the input dataset requirements of our model, the entire field is sliced

469  into 10x10x13 image sets within the 200 West Area. Each image set is assigned a performance

470  ranking for the next 5 years, predicted by the MC3D-CNN model as a potential well location,

471  with a given plume concentration image at any time step, hydrostratigraphic unit, and hydraulic

472  conductivity. Figure 9 presents the resulting performance ranking prediction, which can be used

473  to assist in future well planning.



474

475 Figure 9. 3D performance ranking map. Green, orange, and red indicate the locations of high,

476 medium, and low well performance rankings, respectively. For each category, the level of

477 transparency indicates the confidence of the CNN model prediction. The more transparency, the

478 less confidence there is in CNN model prediction. The transparency of the low-ranking areas

479 (red) is further scaled down to 10% to show the internal high- and medium-ranking areas in the

480 domain.

481 The centroid plume concentration of each grid in Figure 9 was clustered according to its

482 performance ranking. The results are shown in Figure 10 using a violin plot, which is a hybrid of

483 a boxplot and kernel density plot that can visualize the distribution of a numeric variable in

484 different groups. It is not surprising that the low-ranking grids have much lower concentration

485 compared to the medium- and high-ranking grids. The concentration distributions of the medium

486 and high classes are similar. It is interesting to see that many grids with high concentration were

487 ranked as medium while many other grids with lower concentration were ranked as high. This

488 highlights the value of using image-based classification instead of point measurement for

489 decision-making.

490

Figure 10. Violin plot of the cell grid concentration vs. predictive performance ranking (white

dot is the median of each group, the black box is the boxplot with whisker range, the yellow

violin shows the probability density).

494

**5. Conclusion and Discussion**

We propose a DL framework that can predict the future performance of extraction wells by mining the hidden relationship between historical P&T records and contaminant plume distribution, and aquifer hydrogeological properties. Inspired by the success of DL applications in computer image recognition applications, we formulated the well performance prediction as a classification problem similar to image recognition. The resulting DL model can identify key patterns of subsurface properties that control well extraction by learning from pre-labeled historical records and paired subsurface property images. The trained DL model can capture such key patterns around new wells and classify them as high-, medium-, and low-performance-ranking locations. The advantage of this DL method compared to analytical well models is that it can look deep into the heterogeneous nature of aquifer conditions. Although rigorous physical constraints were not explicitly imposed in the DL framework, they were implicitly included by selecting the key physical factors as the training image for model inputs. The deep learning framework is also adaptive to recruit new data whenever they become available to retrain and improve the deep learning model with small computational cost. Compared to numerical F&T models, this approach is much more portable and makes full use of historical P&T records, and thus is suitable for data-driven decision-making and adaptive site management (ASM) to reduce remediation time and cost, as discussed in Section 5.1. The future development of this method is also discussed in Section 5.1, and the limitations are discussed in Section 5.2.

**5.1. Application scenario of the well performance ranking tool**

The data-driven and portable features of the DL approach made it a useful tool for ASM. Here, we refer to ASM as a systematic and iterative management strategy that routinely re-evaluates

518    and prioritizes site remedial actions and characterization activities to expedite the remediation of

519    large and/or complex sites (Demirkanli and Freedman, 2021). Prediction of extraction well

520    performance at specific locations is critical for ASM for evaluating and revising P&T well

521    network design and operation, setting reasonable remedy targets, and estimating remedy costs.

522    Although this method is data-driven, it doesn't incur additional costs for collecting new data

523    because the well mass recovery records are often routinely recorded during the P&T operation,

524    subject to the regulatory requirement. The DL model can also be easily re-trained with better

525    accuracy for continuous planning whenever new data are available. During the routine evaluation

526    of remediation progress, the updated well performance ranking map can assist in planning new

527    well locations, rebalancing the pumping rates for existing wells, or even turning off some low-

528    performance-ranking wells to better use their treating volume for high-performance wells.

529

530    Another application scenario for this performance prediction model is to integrate it with

531    optimization workflows for P&T well network design. There is a long history of developing and

532    applying integrated simulation-based optimization approach for P&T system design (Khan et al.,

533    2004; Maskey et al., 2002; Mayer et al., 2002; McKinney and Lin, 1996; Wagner and Gorelick,

534    1987; Zheng and Wang, 2002). In a typical simulation-based optimization application,

535    optimization search algorithms [e.g., differential evolution (Bayer et al.), generic evolution

536    (Park, 2016), particle swarm (Mategaonkar and Eldho, 2014), firefly (Kazemzadeh-Parsi et al.,

537    2015), and others] are used to drive a groundwater simulator iteratively to check whether the

538    environmental and/or hydraulic constraints were met with certain P&T configuration parameters

539    (e.g., new well location, pumping rates, and pumping duration), and then adjust these

540    configuration parameters accordingly. The computational cost of the optimization problem

31

541  grows exponentially with the number of parameters and makes formal optimization nearly

542  impossible for complex waste sites with a large number of wells. The data-driven MD3D-CNN

543  model is much more portable comparing to expensive numerical F&T models, and it is also more

544  adaptive to recruit new data whenever they become available. These make it ideal as~~The MD3D-~~

545  ~~CNN well performance prediction model can be used as~~ a filtering tool to reduce the number of

546  candidate well locations for the optimization search algorithms so that limited computational

547  resources can be concentrated on more promising well installation plans. The performance

548  prediction model can also be integrated into groundwater simulators (e.g., the Hanford Site's P2R

549  model) and provide direct on-the-fly optimization. In such cases, the model can be used as a

550  wrapper of the groundwater simulator that pauses the simulator periodically and then rebalances

551  the extraction rates among wells based on their performance ranking to achieve better mass

552  recovery.

553  **5.2. Limitations and future development**

554  We demonstrated this performance prediction model using the model-simulated $CCl_4$ of a real

555  complex remediation project located on the Hanford Site. The model simulation results provided

556  a known answer so that accuracy and mismatch of the DL models could be precisely measured

557  and traced. Although the scope of this study is limited to developing and demonstrating the new

558  method, one question to answer is how to apply this method in field as there are no known

559  "subsurface images" such as the exact plume distribution provided by the calibrated P2R model.

560  Our next step will address this question with the following three approaches:

561  •  ~~Geospatial~~ Geostatistical simulation and ensemble prediction. Due to the limited

562     sampling and monitoring data in any real remediation site, it is clearly impossible to have

563     exact plume distribution for the performance prediction model as model input. However,

564     a remediation project always has some type of estimation of the plume distribution,

565     which is essential for decision-making. Such a plume distribution estimation can be used

566     as a training image for the performance ranking tool. An even better method would be to

567     train multiple DL models using the geostatistical realization of the plume distribution

568     (Murray and Bott, 2008) to augment the training pool to correct model bias. The site

569     uncertainty can be reduced by combining all the DL model results to provide a more

570     representative or accurate ensemble prediction.

571   • Incorporating multiple data inputs. Although the demonstration case only used three key

572     aquifer properties as inputs, the MC3D-CNN architecture is flexible and can be readily

573     extended to include other variables. One important potential training image dataset is

574     geophysical investigations, which can provide high spatial-temporal resolution snapshots

575     of subsurface measurements.

576   • Training the DL model with the aid of numerical model results. In sites with very limited

577     subsurface measurement, numerical simulation results can be added as a supplementary

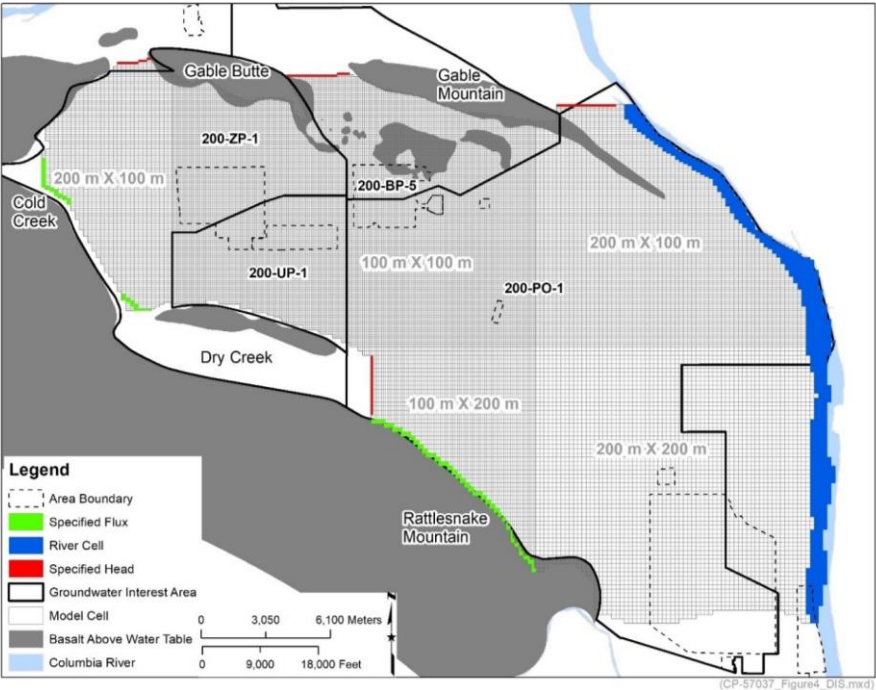578     dataset for model training.

579

580   This study provides a DL framework to make better use of P&T records for future remediation

581   design. P&T remedy monitoring and operational data, along with site investigation information,

582   applied to data-informed approaches such as the one tested in this study, can create opportunities

583   to improve our understanding of contaminant transport, provide flexible tools for site

584   management, streamline decision-making, and potentially reduce remediation costs.
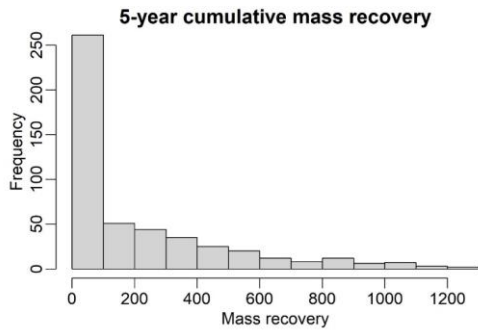
585

586

33

587

**Acknowledgments**

593 **Supporting information**



594

595 Figure SI. P2R version 8.3 model extent and boundary conditions. (Source: CP-57037, Model

596 Package Report: Plateau to River, Version 8.3, Rev. 2.)

597

**5-year cumulative mass recovery**

598

599     Figure S2. Histogram of 5-year cumulative mass recovery.

600     **References**

601     Ameli, A.A., Craig, J.R., 2018. Semi-analytical 3D solution for assessing radial collector well pumping

602          impacts on groundwater–surface water interaction. Hydrology Research, 49(1): 17-26.

603          DOI:10.2166/nh.2017.201

604     Bayer, P., Paly, d.M., Bürger, C.M., 2010. Optimization of high-reliability-based hydrological design

605          problems by robust automatic sampling of critical model realizations. Water Resour. Res., 46(5).

606          DOI:10/d8sv7r

607     Brusseau, M.L., 1996. Evaluation of Simple Methods for Estimating Contaminant Removal by Flushing.

608          Groundwater, 34(1): 19-22. DOI:10.1111/j.1745-6584.1996.tb01860.x

609     Brusseau, M.L., 2013. Use of Historical Pump-and-Treat Data to Enhance Site Characterization and

610          Remediation Performance Assessment. Water Air Soil Pollut, 224(10): 1741. DOI:10/gjhfq9

611     Budge, T., Nichols, W., 2020. Model Package Report: Plateau to River Groundwater Model Version 8.3.

612          CP-57037-Rev.2.

613    Cardiff, M., Liu, X., Kitanidis, P.K., Parker, J., Kim, U., 2010. Cost optimization of DNAPL source and

614        plume remediation under uncertainty using a semi-analytic model. Journal of Contaminant

615        Hydrology, 113(1): 25-43. DOI:10/c75h97

616    Carrera, J., 1993. An overview of uncertainties in modelling groundwater solute transport. Journal of

617        Contaminant Hydrology, 13(1): 23-48. DOI:10.1016/0169-7722(93)90049-X

618    Chen, Y., Liu, G., Huang, X., Meng, Y., 2022. Groundwater Remediation Design Underpinned By

619        Coupling Evolution Algorithm With Deep Belief Network Surrogate. Water Resour Manage,

620        36(7): 2223-2239. DOI:10.1007/s11269-022-03137-w

621    Demirkanli, D.I., Freedman, V.L., 2021. Adaptive Site Management Strategies for the Hanford Central

622        Plateau Groundwater. PNNL-32055.

623    Demirkanli, D.I. et al., 2018. Assessment of Pump-and-Treat System Impacts on 200 West Aquifer

624        Conditions. PNNL--28063, 1490801.

625    EPA, 2002. Groundwater Remedies Selected at Superfund Sites. EPA-542-R-01-022, Washington, D.C.

626    EPA, 2005. Cost-effective design of pump and treat systems. EPA 542-R-05-008, Washington, D.C.

627    Finsterle, S., 2006. Demonstration of optimization techniques for groundwater plume remediation using

628        iTOUGH2. Environmental Modelling & Software, 21(5): 665-680.

629        DOI:10.1016/j.envsoft.2004.11.012

630    Finsterle, S., Zhang, Y., 2011. Solving iTOUGH2 simulation and optimization problems using the PEST

631        protocol. Environmental Modelling & Software, 26(7): 959-968.

632        DOI:10.1016/j.envsoft.2011.02.008

633    Gaur, S., Ch, S., Graillot, D., Chahar, B.R., Kumar, D.N., 2013. Application of Artificial Neural

634        Networks and Particle Swarm Optimization for the Management of Groundwater Resources.

635        Water Resour Manage, 27(3): 927-941. DOI:10/f4mwmv

636 Gaur, S., Chahar, B.R., Graillot, D., 2011. Analytic elements method and particle swarm optimization

637        based simulation–optimization model for groundwater management. Journal of Hydrology,

638        402(3): 217-227. DOI:10/cz25k2

639 Hadley, P.W., Newell, C.J., 2012. Groundwater Remediation: The Next 30 Years. Groundwater, 50(5):

640        669-678. DOI:10.1111/j.1745-6584.2012.00942.x

641 Haley, J.L., Lang, D.J., Herrinton, L., 1989. EPA's approach to evaluating and cleaning up ground water

642        contamination at Superfund sites. Ground Water Monitoring Review; (USA), 9:4.

643        DOI:10.1111/j.1745-6592.1989.tb01027.x

644 Hammond, G.E., Lichtner, P.C., 2010. Field-scale model for the natural attenuation of uranium at the

645        Hanford 300 Area using high-performance computing: MODEL FOR NATURAL

646        ATTENUATION OF URANIUM. Water Resour. Res., 46(9). DOI:10.1029/2009WR008819

647 Hirschmiller, J., Biryukov, A., Groulx, B., Emmerson, B., Quinell, S., 2019. The Importance of

648        Integrating Subsurface Disciplines with Machine Learning when Predicting and Optimizing Well

649        Performance – Case Study from the Spirit River Formation, Day 2 Tue, October 01, 2019,

650        Calgary, Alberta, Canada, pp. D021S025R004. DOI:10/gkq3mb

651 Huang, C., Mayer, A.S., 1997. Pump-and-treat optimization using well locations and pumping rates as

652        decision variables. Water Resour. Res., 33(5): 1001-1012. DOI:10/bjzdwq

653 Kazemzadeh-Parsi, M.J., Daneshmand, F., Ahmadfard, M.A., Adamowski, J., 2015. Optimal

654        Remediation Design of Unconfined Contaminated Aquifers Based on the Finite Element Method

655        and a Modified Firefly Algorithm. Water Resour Manage, 29(8): 2895-2912. DOI:10/f7ccpn

656 Khan, F.I., Husain, T., Hejazi, R., 2004. An overview and analysis of site remediation technologies.

657        Journal of Environmental Management, 71(2): 95-122. DOI:10.1016/j.jenvman.2004.02.003

658 Kontos, Y.N. et al., 2022. Machine learning for groundwater pollution source identification and

659        monitoring network optimization. Neural Comput & Applic. DOI:10.1007/s00521-022-07507-8

660    Kontos, Y.N., Katsifarakis, K.L., 2017. Optimal management of a theoretical coastal aquifer with

661        combined pollution and salinization problems, using genetic algorithms. Energy, 136: 32-44.

662        DOI:10/gb4sv6

663    LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. nature, 521(7553): 436-444.

664    Li, Y., Sun, R., Horne, R., 2019. Deep Learning for Well Data History Analysis. DOI:10.2118/196011-

665        MS

666    Majumder, P., Eldho, T.I., 2016. A New Groundwater Management Model by Coupling Analytic Element

667        Method and Reverse Particle Tracking with Cat Swarm Optimization. Water Resour Manage,

668        30(6): 1953-1972. DOI:10/f8hgcp

669    Majumder, P., Lu, C., 2021. A novel two-step approach for optimal groundwater remediation by coupling

670        extreme learning machine with evolutionary hunting strategy based metaheuristics. Journal of

671        Contaminant Hydrology, 243: 103864. DOI:10.1016/j.jconhyd.2021.103864

672    Maskey, S., Jonoski, A., Solomatine, D.P., 2002. Groundwater Remediation Strategy Using Global

673        Optimization Algorithms. Journal of Water Resources Planning and Management, 128(6): 431-

674        440. DOI:10.1061/(ASCE)0733-9496(2002)128:6(431)

675    Mategaonkar, M., Eldho, T.I., 2014. Multiobjective Groundwater Remediation Design Using a Coupled

676        MFree Point Collocation Method and Particle Swarm Optimization. Journal of Hydrologic

677        Engineering, 19(6): 1259-1263. DOI:10/f54mkm

678    Matott, L.S., Rabideau, A.J., Craig, J.R., 2006. Pump-and-treat optimization using analytic element

679        method flow models. Advances in Water Resources, 29(5): 760-775. DOI:10/fq9dkh

680    Mayer, A.S., Kelley, C.T., Miller, C.T., 2002. Optimal design for problems involving flow and transport

681        phenomena in saturated subsurface systems. Advances in Water Resources, 25(8): 1233-1256.

682        DOI:10.1016/S0309-1708(02)00054-4

683    Mayer, K.U., Blowes, D.W., Frind, E.O., 2001. Reactive transport modeling of an in situ reactive barrier

684          for the treatment of hexavalent chromium and trichloroethylene in groundwater. Water Resour.

685          Res., 37(12): 3091-3103. DOI:10.1029/2001WR000234

686    McConnell, L. et al., 2022. Forecasting Groundwater Contaminant Plume Development Using Statistical

687          and Machine Learning Methods. Groundwater Monit R. DOI:10.1111/gwmr.12523

688    McKinney, D.C., Lin, M.-D., 1996. Pump-and-Treat Ground-Water Remediation System Optimization.

689          Journal of Water Resources Planning and Management, 122(2): 128-136.

690          DOI:10.1061/(ASCE)0733-9496(1996)122:2(128)

691    McMahon, A., Heathcote, J., Carey, M., Erskine, A., 2001. Guide to good practice for the development of

692          conceptual models and the selection and application of mathematical models of contaminant

693          transport processes in the subsurface. National Groundwater & Contaminated Land Centre.

694          Environment Agency. UK. Report NC/99/38, 2.

695    Meray, A.O. et al., 2022. PyLEnM: A Machine Learning Framework for Long-Term Groundwater

696          Contamination Monitoring Strategies. Environ. Sci. Technol., 56(9): 5973-5983.

697          DOI:10.1021/acs.est.1c07440

698    Minsker, B., Zhang, Y., Greenwald, R., Peralta, R., Zheng, C., 2004. Application of Flow and Transport

699          Optimization Codes to Groundwater Pump and Treat Systems- Volume III, Fort Belvoir, VA.

700    Mo, S., Zabaras, N., Shi, X., Wu, J., 2019. Deep Autoregressive Neural Networks for High-Dimensional

701          Inverse Problems in Groundwater Contaminant Source Identification. Water Resour. Res., 55(5):

702          3856-3881. DOI:10.1029/2018WR024638

703    Murray, C., Bott, Y.-J., 2008. Revised Geostatistical Analysis of the Inventory of Carbon Tetrachloride in

704          the Unconfined Aquifer in the 200 West Area of the Hanford Site. DOI:10.2172/945229

705    National Research, C., 1994. Alternatives for Ground Water Cleanup.

706 National Research, C., 2013. Alternatives for Managing the Nation's Complex Contaminated

707          Groundwater Sites.

708 Neville, C., Tonkin, M., 2004. Modeling multiaquifer wells with MODFLOW. Ground water, 42: 910-9.

709          DOI:10.1111/j.1745-6584.2004.t01-9-.x

710 Park, Y.-C., 2016. Cost-effective optimal design of a pump-and-treat system for remediating groundwater

711          contaminant at an industrial complex. Geosci J, 20(6): 891-901. DOI:10/gjhfmk

712 Rao, C., Liu, Y., 2020. Three-dimensional convolutional neural network (3D-CNN) for heterogeneous

713          material homogenization. Computational Materials Science, 184: 109850.

714          DOI:10.1016/j.commatsci.2020.109850

715 Razavi, S., Tolson, B.A., Burn, D.H., 2012. Review of surrogate modeling in water resources. Water

716          Resour. Res., 48(7). DOI:10.1029/2011WR011527

717 Ren, H., Cromwell, E., Kravitz, B., Chen, X., 2022. Technical note: Using long short-term memory

718          models to fill data gaps in hydrological monitoring networks. Hydrology and Earth System

719          Sciences, 26(7): 1727-1743. DOI:10.5194/hess-26-1727-2022

720 Rodriguez-Galiano, V., Mendes, M.P., Garcia-Soldado, M.J., Chica-Olmo, M., Ribeiro, L., 2014.

721          Predictive modeling of groundwater nitrate pollution using Random Forest and multisource

722          variables related to intrinsic and specific vulnerability: A case study in an agricultural setting

723          (Southern Spain). Science of The Total Environment, 476-477: 189-206.

724          DOI:10.1016/j.scitotenv.2014.01.001

725 Rogers, L.L., Dowla, F.U., 1994. Optimization of groundwater remediation using artificial neural

726          networks with parallel solute transport modeling. Water Resour. Res., 30(2): 457-481.

727          DOI:10.1029/93WR01494

728 Shen, C., 2018. A Transdisciplinary Review of Deep Learning Research and Its Relevance for Water

729          Resources Scientists. Water Resour. Res., 54(11): 8558-8593. DOI:10.1029/2018WR022643

730    Simonyan, K., Zisserman, A., 2015. Very Deep Convolutional Networks for Large-Scale Image

731        Recognition. arXiv:1409.1556 [cs].

732    Singh, A., Minsker, B.S., 2008. Uncertainty-based multiobjective optimization of groundwater

733        remediation design. Water Resour. Res., 44(2). DOI:10.1029/2005WR004436

734    Soriano, M.A. et al., 2021. Assessment of groundwater well vulnerability to contamination through

735        physics-informed machine learning. Environ. Res. Lett., 16(8): 084013. DOI:10.1088/1748-

736        9326/ac10e0

737    Steefel, C.I. et al., 2015. Reactive transport codes for subsurface environmental simulation. Comput

738        Geosci, 19(3): 445-478. DOI:10.1007/s10596-014-9443-x

739    Sun, A.Y., 2018. Discovering State-Parameter Mappings in Subsurface Models Using Generative

740        Adversarial Networks. Geophysical Research Letters, 45(20): 11,137-11,146.

741        DOI:10.1029/2018GL080404

742    Sváb, M., Zilka, M., Müllerová, M., Kocí, V., Müller, V., 2008. Semi-empirical approach to modeling of

743        soil flushing: model development, application to soil polluted by zinc and copper. Sci Total

744        Environ, 392(2-3): 187-197. DOI:10.1016/j.scitotenv.2007.12.001

745    Tahmasebi, P., Kamrava, S., Bai, T., Sahimi, M., 2020. Machine learning in geo- and environmental

746        sciences: From small to large scale. Advances in Water Resources, 142: 103619.

747        DOI:10.1016/j.advwatres.2020.103619

748    Tartakovsky, A.M., Marrero, C.O., Perdikaris, P., Tartakovsky, G.D., Barajas-Solano, D., 2020. Physics-

749        Informed Deep Neural Networks for Learning Parameters and Constitutive Relationships in

750        Subsurface Flow Problems. Water Resour. Res., 56(5): e2019WR026731.

751        DOI:10.1029/2019WR026731

752    Truex, M. et al., 2017. Performance Assessment of Pump-and-Treat Systems. Groundwater Monit R,

753        37(3): 28-44. DOI:10/gc2hr8

754    Tsang, C.-F., Neretnieks, I., Tsang, Y., 2015. Hydrologic issues associated with nuclear waste

755        repositories. Water Resour. Res., 51(9): 6923-6972. DOI:10.1002/2015WR017641

756    Wagner, B.J., Gorelick, S.M., 1987. Optimal groundwater quality management under parameter

757        uncertainty. Water Resour. Res., 23(7): 1162-1174. DOI:10.1029/WR023i007p01162

758    Wang, N., Chang, H., Zhang, D., 2021. Deep-Learning-Based Inverse Modeling Approaches: A

759        Subsurface Flow Example. Journal of Geophysical Research: Solid Earth, 126(2):

760        e2020JB020549. DOI:10.1029/2020JB020549

761    White, M.D., Oostrom, M., 2003. STOMP subsurface transport over multiple phases version 3.0 User's

762        guide.

763    Wu, C., Fang, C., Wu, X., Zhu, G., 2020. Health-Risk Assessment of Arsenic and Groundwater Quality

764        Classification Using Random Forest in the Yanchi Region of Northwest China. Expo Health,

765        12(4): 761-774. DOI:10.1007/s12403-019-00335-7

766    Wu, J., Zeng, X., 2013. Review of the uncertainty analysis of groundwater numerical simulation. Chin.

767        Sci. Bull., 58(25): 3044-3052. DOI:10.1007/s11434-013-5950-8

768    Yadav, B., Mathur, S., Ch, S., Yadav, B.K., 2018. Data-based modelling approach for variable density

769        flow and solute transport simulation in a coastal aquifer. Hydrological Sciences Journal, 63(2):

770        210-226. DOI:10.1080/02626667.2017.1413491

771    Yan, S., Minsker, B., 2006. Optimal groundwater remediation design using an Adaptive Neural Network

772        Genetic Algorithm. Water Resour. Res., 42(5). DOI:10.1029/2005WR004303

773    Yin, J., Tsai, F.T.C., 2020. Bayesian set pair analysis and machine learning based ensemble surrogates for

774        optimal multi-aquifer system remediation design. Journal of Hydrology, 580: 124280.

775        DOI:10/gk7s8f

776    Yu, X. et al., 2020. Deep learning emulators for groundwater contaminant transport modelling. Journal of

777        Hydrology, 590: 125351. DOI:10.1016/j.jhydrol.2020.125351

778    Zhao, X. et al., 2019. A Multi-Branch 3D Convolutional Neural Network for EEG-Based Motor Imagery

779        Classification. IEEE Transactions on Neural Systems and Rehabilitation Engineering, 27(10):

780        2164-2177. DOI:10.1109/TNSRE.2019.2938295

781    Zheng, C., Wang, P.P., 2002. A Field Demonstration of the Simulation Optimization Approach for

782        Remediation System Design. Groundwater, 40(3): 258-266. DOI:10.1111/j.1745-

783        6584.2002.tb02653.x

784    Zounemat-Kermani, M., Batelaan, O., Fadaee, M., Hinkelmann, R., 2021. Ensemble machine learning

785        paradigms in hydrology: A review. Journal of Hydrology, 598: 126266.

786        DOI:10.1016/j.jhydrol.2021.126266

787