# Journal Pre-proofs

Valdrich Jude Fernandes, Perry de Louw, Ruud Bartholomeus, Coen Ritsema

Please cite this article as: Jude Fernandes, V., de Louw, P., Bartholomeus, R., Ritsema, C., Machine learning for faster estimates of groundwater response to artificial aquifer recharge, *Journal of Hydrology* (2024), doi: https://doi.org/10.1016/j.jhydrol.2024.131418

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Machine learning for faster estimates of groundwater response to artificial aquifer recharge

Valdrich Jude Fernandes[1*], Perry de Louw[2,1], Ruud Bartholomeus[3,1] and Coen Ritsema[1]

1. Soil Physics and Land Management, Wageningen University & Research, Wageningen, the Netherlands
2. Deltares, Utrecht, the Netherlands
3. KWR Water Research Institute, Nieuwegein, the Netherlands

* Corresponding author (email address: valdrich.fernandes@wur.nl, postal address: P.O. Box 47, 6700 AA Wageningen, The Netherlands)

## Abstract

Groundwater models are a valuable tool in optimising the decisions influencing groundwater flow. Spatially distributed models represent the groundwater level in the entire area from where essential information can be extracted, directly aiding in the decision-making process. However, these models are time-consuming, limiting the number of scenarios that can be considered. This study explores different machine learning (ML) models as faster alternatives to predict the increase in steady-state groundwater head due to artificial recharge in the unconfined aquifer while considering a wider spatial extent (832 columns x 1472 rows totalling 765 km$^2$) than previous ML groundwater models. We trained three ML models (encoder-decoder, U-Net, and attention U-Net) with various hypothetical artificial recharge sites (100, 300, 500, and 1000 sites) in the Baakse Beek catchment (the Netherlands), using a detailed numerical groundwater model, AMIGO. The applied recharge rate along with geo-hydrological properties from the AMIGO baseline run were used as inputs to the ML models. The properties' permutation importance indicated that all properties of the first aquifer were important to predicting the response and were included when training the ML models. All three ML models improved with additional training sites but showed limited benefits from more than 500 recharge sites. Of the three ML models, U-Net and Attention U-Net outperformed the encoder-decoder. These two models achieved Nash-Sutcliffe efficiency (NSE) of more than 0.8 when trained with 300 or more recharge sites. U-Net trained on 1000 recharge sites had the highest overall NSE of 0.95. U-Net better captures input features with highly variable spatial characteristics, such as rivers and drains which influence the maximum height of the groundwater response. The model captured the influence of the input features on the response, reproducing the response patterns across the entire catchment. Finally, we showed that the trained ML models are faster than the numerical model, predicting within 0.24 seconds (97th percentile), making it ideal for optimising decisions.

50

## 1. Introduction

52  In the context of international policy frameworks like the European
53  Water Framework Directive and Natura 2000, water management
54  authorities have multiple targets they need to meet. The droughts
55  of 2018-2020, which set a new benchmark in Europe (Rakovec et al.,
56  2022), increased the urgency to take appropriate measures
57  (Bartholomeus et al., 2023). Although the events are considered
58  rare in the current climate, future climate change could exacerbate
59  such events (Aalbers et al., 2023; Balting et al., 2021; Lehner et al.,
60  2017; Pronk et al., 2021; van der Wiel et al., 2021). Even in deltas
61  like the Netherlands droughts cause serious risks for nature,
62  agriculture, infrastructure and drinking water availability, which
63  resulted in drought-related policy actions like "Water and soil
64  leading in land use planning" (Bartholomeus et al., 2023). One of the
65  reasons for this vulnerability is the expansion of the surface
66  drainage network and the increased exploitation of groundwater
67  resources (Ahmadalipour et al., 2019; Bartholomeus et al., 2023;
68  Castle et al., 2014; de Wit et al., 2022; Thatch et al., 2020; Thomas
69  and Famiglietti, 2019; Witte et al., 2018).

70  The Pleistocene uplands of the Netherlands have recently faced
71  severe rainfall deficits (Brakkee et al., 2022; Philip et al., 2020),
72  increasing the reliance on surface and groundwater for irrigation.
73  This has increased the strain on the limited water available for
74  nature (van den Eertwegh et al., 2020). Long-term structural
75  changes are identified to be more effective at reducing the strain
76  than reactive, ad-hoc remedies during droughts. Van den Eertwegh
77  et al. 2020 recommend increasing freshwater availability through
78  more sustainable drainage networks, reducing groundwater
79  abstraction, and increasing groundwater recharge.

80  Managed aquifer recharge (MAR) can increase freshwater
81  availability during dryer periods by storing water surplus from the
82  wetter periods in the subsurface (Dillon et al., 2020, 2019; Hartog
83  and Stuyfzand, 2017). It is often categorised into infiltration, direct
84  injection, and filtration techniques (Casanova et al., 2016); we focus
85  on infiltration techniques that recharge the water table from
86  infiltration basins or subsurface infiltration systems, often making
87  them the cheapest technique. However, water managers need to
88  identify the optimal location, recharge rate and combination of the
89  recharge sites when designing the solution which is often done using
90  a numerical groundwater model. These models use a set of
91  mathematical equations to estimate the flow of water within a grid
92  that represents the hydrological system by their characteristics, such

2

93    as the aquifer's transmissivity, resistance and the surface drainage
94    network. However, they are complex and simulating multiple
95    scenarios for optimisation is time-consuming, limiting structured
96    exploration and selection of potential recharge sites across an area.
97    To facilitate the exploration of suitable recharge sites, there is a
98    need for fast calculating tools to estimate the effect of managed
99    aquifer recharge quickly. For such optimisation applications, a less
100   accurate but faster option with interpretable results could be more
101   suitable (Newman, 1996).

102   Such an option could be a surrogate model, which is a simplified
103   representation of a complex, higher-order model (Wang et al.,
104   2014). Reduced order models have been applied in groundwater
105   modelling as surrogate models for their computational efficiency
106   (Boyce et al., 2015; Dey and Dhar, 2020; Stanko et al., 2016;
107   Vermeulen et al., 2004). Proper orthogonal decomposition, a
108   common reduced-order modelling method, identifies the lower
109   dimensional basis that captures the high-dimensional dynamics of
110   the system. Vermeulen et al. (2004) have demonstrated its
111   applicability in reproducing groundwater heads in a linear system. In
112   a realistic case study, they achieved a relative mean absolute error
113   of less than 6% while realising a 625x speed up. However, these
114   attempts have been made for confined conditions with linear
115   behaviour. Boyce et al. (2015) and Stanko et al. (2016) expanded
116   this technique to unconfined aquifers, increasing the nonlinear
117   behaviour due to the boundary conditions such as rivers. While
118   more realistic, they are still limited to small synthetic systems with
119   less than 200 by 200 cells. Furthermore, proper orthogonal
120   decomposition models are limited to the location used to calculate
121   the reduced space.

122   Machine learning (ML) has recently been a frequently used
123   surrogate model as a universal function approximator. It can learn
124   nonlinear relations in the data, which can be the results from
125   existing numerical models. It has been used to reproduce models in
126   fluid dynamics (Brunton et al., 2020), material science
127   (Papadopoulos et al., 2018) and earth system models (Kim et al.,
128   2015; Weber et al., 2019), among others. Deep learning models
129   have been used in groundwater modelling to forecast the head at
130   wells (Malik and Bhagwat, 2021; Müller et al., 2021; Tao et al.,
131   2022). Asher et al. (2015) and Miro et al. (2021) recognised the lack
132   of spatially distributed representation of groundwater surrogates.
133   Since then, some authors have demonstrated the applicability of the
134   convolutional encoder-decoder model, which satisfies this requisite
135   (He et al., 2021; Mo et al., 2019; Taccari et al., 2022). However,
136   these applications are also limited to small synthetic systems.

137   Artificial groundwater recharge affects the groundwater head in a
138   large spatial area. This entire spatial extent needs to be captured by
139   the ML model. The applicability of the above ML models at
140   reproducing the results from a numerical groundwater model with
141   actual subsurface properties of an aquifer has not been

3

142 demonstrated yet. Furthermore, the ML model can be more
143 specialised and represent the priorities of the optimisation
144 challenge rather than a model reproducing all details of the system.
145 We investigate the performance of three ML models for a
146 catchment within the sandy uplands of the Netherlands and
147 quantify the effect of artificial recharge for all possible locations
148 within the area. The ML models' output is the increase in the steady-
149 state phreatic groundwater head, henceforth groundwater
150 response, to applied recharge sites in the Baakse Beek catchment in
151 the Netherlands. The hydrological properties and the results from a
152 detailed numerical model (AMIGO) are used to train the ML models.
153 The ML model is trained on the geo-hydrological properties of the
154 first aquifer for a wider domain size of 1472 columns by 832 rows at
155 a 25x25 m resolution representing a 765 km$^2$ area. In doing so, we
156 consider various combinations of geo-hydrological properties within
157 the catchment and their impact on the performance of the
158 surrogate model at predicting the steady-state groundwater head
159 response to artificial recharge. These steps are further elaborated in
160 the methodology and through the flow chart in Figure 1. The central
161 questions this study aims to answer are:

162     1.  Is the surrogate model able to reproduce the steady-state
163         groundwater head response to artificial recharge with
164         sufficient accuracy?
165     2.  Which physical characteristics are required to capture the
166         steady-state response of the groundwater head to artificial
167         recharge in a surrogate model trained on the results of a
168         numerical model?
169     3.  How much training data is needed to train the surrogate
170         model to sufficient accuracy?

171 In addressing these questions, this paper aims to aid future
172 modellers in designing more accurate ML models for scenario
173 optimisations. These questions remain relevant even through the
174 fast advancement in artificial intelligence and ML. Multiple geo-
175 hydrological properties represent the subsurface, but identifying the
176 most relevant properties could help the ML model capture the
177 relation between them and reduce overfitting. Furthermore, we
178 want to minimize the number of slow numerical model runs. This
179 paper compares the performance of the ML model when trained on
180 datasets of various sizes. This offers an estimate of the number of
181 scenarios needed to train the ML models and the effect of additional
182 scenarios on the predicted groundwater response. Comparing three
183 ML models with increasing complexity offers a more general view of
184 answering the above questions and model complexity necessary to
185 represent the relation between the recharge rate, hydrogeological
186 properties and the groundwater response.
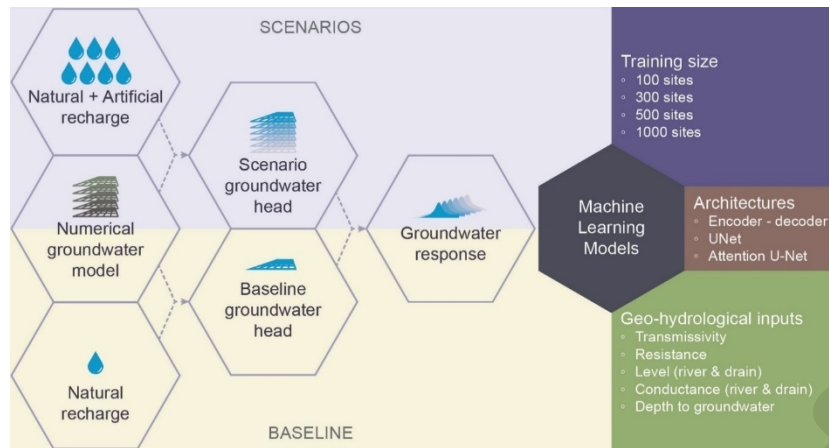
187 ## 2. Methodology

188 The research methodology consists of two main parts: numerical
189 modelling and machine learning modelling (Figure 1). The goal is to

4

190     use the numerical model to simulate a baseline steady-state
191     scenario of natural recharge and steady-state scenarios with
192     artificial recharge at sites across the study catchment. The
193     difference in the groundwater heads between the artificial recharge
194     scenarios and the baseline scenario is the groundwater response to
195     the artificial recharge. A machine learning model is trained to
196     reproduce this response. The rate of artificial recharge (5-25
197     mm/day) and site size (0.01-1 $km^2$) are selected randomly using
198     Latin Hypercube Sampling to represent the entire range of potential
199     recharge sites. Orthogonal Array-based Latin Hypercube Sampling is
200     used to select the site location, within the model extent, as it
201     samples the location more uniformly.

202     The ML models are trained to reproduce the steady-state
203     groundwater head response due to artificial recharge from the
204     numerical groundwater model, AMIGO. These ML models are
205     trained on training datasets of various artificial recharge
206     realizations. Each realization contains six inputs from the AMIGO
207     baseline run: (1) the artificial recharge rate, (2) baseline
208     groundwater depth, (3) river stage and drain level relative to the
209     baseline groundwater head, (4) river conductivity, (5) transmissivity
210     of the first aquifer and (6) hydraulic resistance below the aquifer.
211     The inputs were included based on their permutation importance in
212     estimating three key characteristics of the steady-state groundwater
213     head response to artificial recharge, namely the maximum, area,
214     and total response. The ML model performance is also assessed on
215     the same three key characteristics as they describe the most
216     relevant properties of the response to optimize.

217     For steady-state simulations, the storage coefficient is zero by
218     definition, thus not an input of the numerical model simulations,
219     and therefore also not included in the inputs for the ML model. It
220     should be noted, however, that in transient simulations the storage
221     coefficient will be another system characteristic that importantly
222     influences aquifer storage capacity to artificial recharge.
223     Additionally, using the storage coefficient from a transient model
224     lets us estimate the extra volume of water which can be stored by
225     the artificial recharge, based on the simulated head differences.

226     Three ML models are trained using the listed inputs: encoder-
227     decoder, U-Net and Attention U-Net, with increasing numbers of
228     recharge sites: 100, 300, 500 and 1000. These models are designed
229     to be increasingly complex, with the Attention U-Net having the
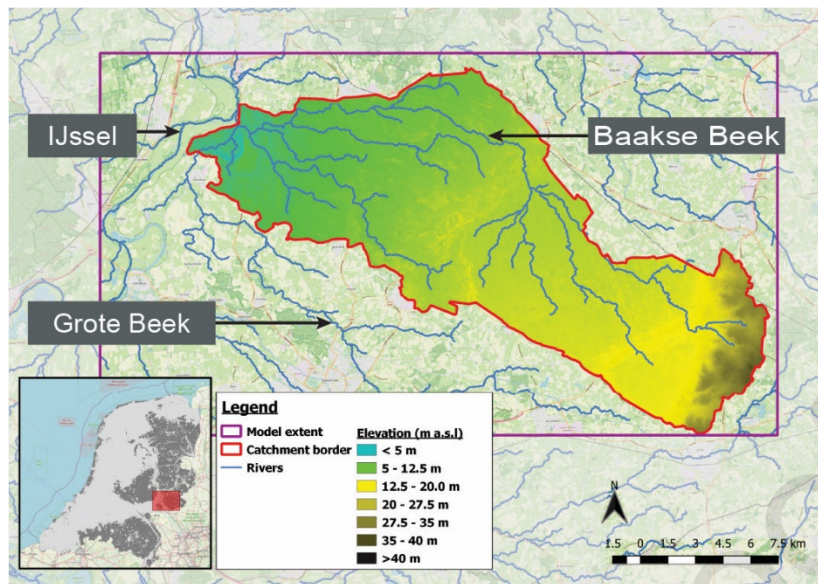230     highest number of parameters.

231

232 *Figure 1 Steps performed before training the machine learning models to reproduce*
233 *the groundwater response to additional aquifer recharge. The groundwater*
234 *response is the increase in the steady-state groundwater head in the scenarios with*
235 *the artificial recharge over the baseline scenario. The scenarios were simulated*
236 *using the numerical groundwater model AMIGO. We compared the importance of*
237 *different geo-hydrological inputs, machine learning model architectures and the*
238 *number of scenarios necessary to train the model.*

## 2.1. Numerical Modelling

240 The ML model is designed to reproduce the steady-state response
241 to additional artificial recharge in the Baakse Beek Catchment east
242 of the Netherlands, as simulated by the AMIGO numerical
243 groundwater model. The catchment drains an area of 262.5 km$^2$ into
244 the IJssel, a distributary of the river Rhine (Figure 2). This catchment
245 is in the Netherlands' higher sandy region, insert in Figure 2,
246 characterized by a 200m-thick sequence of Pleistocene sands
247 intercalated with thin clay beds, which become thicker towards the
248 west. It is mainly composed of coarse-textured glacial and beach
249 deposits (Hijma, 2017; Sevink and Koopman, 2020), which are highly
250 transmissive.

251

*Figure 2 Map of the study area, Baakse Beek catchment, in the sandy region of*
*eastern Netherlands (the dark grey region in insert)*

252
253

254  The Baakse Beek catchment is represented in the spatially
255  distributed regional groundwater model AMIGO (Actueel Model
256  Instrument Gelderland Oost v3.1) which covers the eastern region of
257  the province of Gelderland. It is widely used by the regional water
258  management authority Rijn en IJssel, province of Gelderland,
259  drinking water companies, and consultancies. The model was
260  calibrated and validated by its maintainers (Vreugdenhil, 2021).
261  Within Baakse Beek, the maintainers determined the modelled
262  average low groundwater level is 5 cm higher and the average high
263  groundwater level is 22 cm lower than the observed levels in 2008-
264  2016.

265  The AMIGO model consists of 15 layers, represented by their
266  transmissivity and the hydraulic resistance between them at a 25m
267  resolution. The hydraulic resistance is calculated as saturated
268  thickness divided by the vertical hydraulic conductivity of the
269  aquifers and the resistive layer between them. This model, which
270  includes tile drainage, ditches, streams, and extraction wells, is
271  implemented in iMOD (Vermeulen et al., 2021) for MODFLOW-2005
272  (Harbaugh, 2005). In AMIGO, installed tile drainage are modelled
273  using the DRN package in MODFLOW while ditches and streams are
274  modelled with the RIV package. These two packages together
275  represent the surface water network that drains the groundwater.
276  To help the model capture the effect of the surface water network,
277  the two packages are combined in a common input, referred to here
278  as DRN and RIV. The AMIGO model was then cropped to a rectangle
279  containing the Baakse Beek catchment. A fixed head boundary
280  condition was defined along the edge of a rectangle surrounding the
281  study catchment. The boundary is maintained at a distance of three
282  times the leakage factor from the catchment's boundary to ensure
283  that the boundary does not significantly influence the calculated
284  response to recharge sites within the catchment. The groundwater

285 head from a steady-state run with long-term temporal average
286 natural recharge is used as the model's initial and boundary
287 conditions.

## 2.2. Machine learning models

289 *General modelling task* - The results from the numerical model
290 scenarios are used to train a machine learning (ML) model using a
291 surrogate modelling approach. In this approach, a surrogate model (
292 $f'$) is used to approximate the results ($y$) of a complex model ($f$) by
293 reproducing its outputs. However, in this study, the ML model
294 predicts the difference between the results from a natural recharge
295 scenario ($f_o$) and the artificial recharge scenario ($f_s$) (equation 2)
296 rather than the complex model results directly (equation 1). This
297 increases the relevance of the surrogate model to the scenario
298 optimization task. Predicting the difference also reduces the output
299 range, improving the training process for ML models.

300
$$f'(x^{N_x^* \times H \times W}) \approx f(x^{N_x \times H \times W}) = y^{N_y \times H \times W} \tag{1}$$

301
$$f'(x^{N_x^* \times H \times W}) \approx f_s(x^{N_x \times H \times W}) - f_o(x^{N_x \times H \times W}) = y^{N_y \times H \times W} \tag{2}$$

302 The spatially distributed models, like the AMIGO model, use $N_x$
303 geohydrological features of size $H \times W$ to predict $N_y$ outputs. The
304 ML model aims to estimate the same results based on fewer input
305 features ($N_x^*$) than the numerical model. This reduction in input
306 features helps train a more generalised and representative ML
307 model (Kutz and Brunton, 2022). However, the model needs
308 minimum input features to capture all relevant relations. The
309 numerical groundwater model requires 105 two dimensional
310 features, while the ML models reproduce the response based on 6
311 input features.

312 *Model architecture* - Convolutional neural networks (CNN) (LeCun et
313 al., 2015; Lecun et al., 1998), a popular ML model for image
314 processing, are utilised in this study. These networks are especially
315 suited for learning the local relations within the input features,
316 which can influence the groundwater system in neighbouring grids.
317 In the context of this paper, a feature is a measurable property that
318 is input to the subsequent model layers. CNNs combine multiple
319 layers to extract different features from the input, using trainable
320 weight matrices (filters) that consider the surrounding cells of the
321 cell of interest. Deeper layers in CNNs extract higher-order features,
322 while initial layers extract elementary features. These higher-order
323 features are crucial to capture interactions between the input
324 features (Lerman et al., 2021). In addition to the layers with filters,
325 CNNs also consist of convolutional, upsampling, batch normalisation
326 (Ioffe and Szegedy, 2015), leaky ReLU (Maas et al., 2013) and
327 dropout layers (Srivastava et al., 2014), which together enable
328 learning nonlinear relations between the input features.

8

329     This study compares three ML models: encoder-decoder, U-NET,
330     and attention U-NET. The three models are based on an encoder-
331     decoder architecture. This architecture consists of encoder blocks
332     (left block in Figure 3B) that learn the context in the input features
333     and decoder blocks (right block in Figure 3B) that reconstruct the
334     results from the learned context. The models differ in their encoder-
335     decoder architecture, with variations in the number of filters.

336     The three encoder-decoder models share the same set of 6 input
337     features. The inputs are two-dimensional matrices, i.e. spatially
338     distributed values, of artificial recharge rate, aquifer transmissivity,
339     vertical hydraulic resistance, DRN and RIV conductance, DRN and
340     RIV stage relative to the groundwater head of the baseline run and
341     the depth to the groundwater head of the baseline run (Figure 3A).
342     The features are selected to represent the groundwater flow within
343     the phreatic aquifer, whose importance is confirmed based on
344     permutation importance. These features are passed to the first
345     down sampling block, which generates 32 features. The number of
346     features is doubled by subsequent down sampling blocks, up to 128
347     features. This limit was set to reduce the memory requirements for
348     training the models. After the encoder block, a bottleneck (bottom
349     Figure 3B) containing two convolution layers with 256 features was
350     added, which improves the extent to where the recharge site
351     influences the response.

352     Following the bottleneck, five decoder blocks are used to
353     reconstruct the output with decreasing numbers of features (128,
354     128, 128, 64, and 32) in reverse order compared to the encoder
355     blocks. The final up-sampling to the input dimensions was done
356     using a convolution transpose and a convolution layer. The
357     convolution transpose consists of 8 filters of size 4x4 with stride 2,
358     while the convolution layer produces one feature with a 1x1 filter.
359     Finally, a leaky ReLU activation function is applied to scale back
360     negative values and better represent the output.

361     *Encoder* - The encoder block learns context with five down sampling
362     blocks (left half of Figure 3B). Each block reduces the input's height
363     and width by half using convolutional layers of 5x5 filters and a
364     stride of two and zero padding. These layers are followed by batch
365     normalisation, leaky ReLU activation, and a dropout rate of 10%. The
366     batch normalisation layer normalises the features with a mean of 0
367     and a unit standard deviation. The dropout layer replaces a random
368     subset of the features with 0 during each iteration of the training
369     process, hiding those features and reducing overfitting. The leaky
370     ReLU activation introduces non-linearity to the model by scaling
371     negative values with a slope of 0.2. It is preferred over ReLU, which
372     only considers positive values, to avoid the 'dying ReLU problem'
373     due to which the model weights do not update through gradient
374     descent. The learned features in the encoder are then passed to the
375     decoder, which recreates the response based on the learned
376     context.

9

377 *Decoder* - The decoder block increases the dimension of the features
378 back to that of the input through five upsampling blocks (right half
379 of Figure 3B). The three models differ in their decoders. The
380 simplest of the models is the encoder-decoder, where each
381 upsampling block consists of a convolutional layer followed by
382 bilinear upsampling, leaky ReLU activation (slope 0.2), batch
383 normalisation, and dropout (rate 10%). The convolutional layer uses
384 5x5 filters, a stride of 1, and zero padding.

385 U-Net trains on higher-level features directly from the encoder and
386 context from the deepest part of the network through skip
387 connections. These connections join the feature from the encoder
388 with upsampled features from deeper parts of the network. The
389 combined features are then processed by convolutional layers,
390 batch normalisation, leaky ReLU activation, and dropout layers like
391 in the encoder-decoder.

392 The upsampling blocks in Attention U-Net (Oktay et al., 2018) are
393 similar to that in U-NET. However, it learns to focus on specific
394 regions in the higher-level features using an attention block.
395 Information is extracted from the two sources of features,
396 upsampled contextual features and the higher-level features, using
397 convolutional layers with 3x3 filters, a stride of 1 and zero padding.
398 Additive importance is then calculated based on the information
399 learnt from the two features, and non-linearity is added to the
400 importance with ReLU activation. From these, a single importance
401 weightage is calculated using a convolution layer with a 1x1 filter
402 and stride one and sigmoid activation, which scales the importance
403 between 0 and 1. The detailed features are multiplied with
404 corresponding weights to enhance the relevance of important
405 regions and they are then concatenated with the upsampled
406 contextual feature. This concatenated feature is then passed
407 through the convolutional layers, leaky ReLU activation, batch
408 normalisation, and dropout layers, similar to the previous models.
409 Note that the attention U-Net has half the number of filters as the
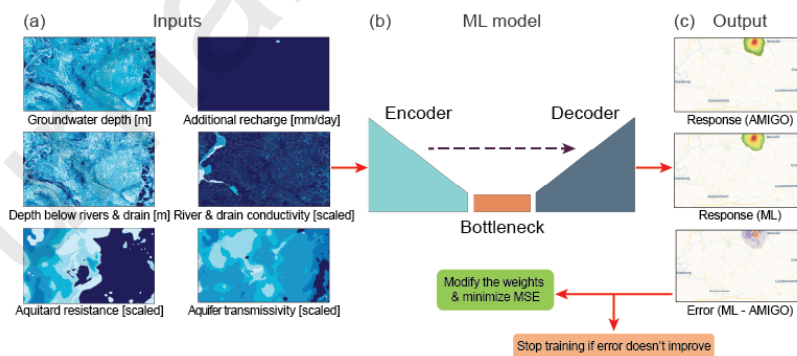410 other two models to stay within the memory limits.

411 *Custom loss function* - The models utilised in this study employ a
412 type of machine learning called supervised learning (Bishop, 2006),
413 which aims to learn a mapping between inputs and outputs based
414 on labelled examples. Specifically, the ML model outputs are
415 compared to groundwater response from the numerical model,
416 AMIGO. The model's performance is evaluated using a loss function
417 such as mean squared error (MSE) to update the model parameters
418 through gradient descent. The training procedure is monitored by
419 tracking the ML model's performance on a validation dataset, which
420 is concluded when the loss does not improve through the training
421 iterations. This validation dataset consists of the AMIGO simulation
422 results from 100 recharge sites that the ML model was not trained
423 on.

424 $$\mathcal{L}_{MSE} = \frac{1}{N^*}\Sigma(y - \hat{y})^2 \tag{3}$$

425 $$\mathcal{L} = (1 - \alpha) * MSE_0 + (\alpha) * MSE_r \tag{4}$$

426 The loss function was modified to help make it more suitable for the
427 task. The target variable from AMIGO is sparse, consisting of
428 multiple cells with no groundwater response, which can lead to the
429 model primarily predicting zeros. The dying ReLU problem (Lu et al.,
430 2020) further exacerbates this problem. To address this, the mean
431 squared error loss ($\mathcal{L}_{MSE}$ in Eq 3) was split into two components in
432 Eq 4: MSE for predictions where there is no response to the applied
433 artificial recharge ($MSE_0$), and MSE for predictions of the response (
434 $MSE_r$). In Eq 3, $y$ and $\hat{y}$ are the groundwater response from AMIGO
435 and the ML model respectively, while $N^*$ is the number of input sites
436 in the iteration. The final loss ($\mathcal{L}$) is a weighted sum of these two
437 components, controlled by a hyperparameter $\alpha$ (Eq 4). This loss
438 function offers several advantages: it balances the error between
439 the overrepresented zeros and the response, unlike a mask, it is still
440 sensitive to predictions away from the site, and the tuneable
441 parameter $\alpha$ allows the relative importance of the two components
442 to be adjusted to reflect the priorities of the use case.

443 Based on this loss, the ML model parameters were iteratively
444 updated using the ADAM optimiser (Kingma and Ba, 2014), with a
445 learning rate schedule. Each iteration consisted of 8 recharge sites
446 (batch size = 8). The initial learning rate was set to 0.002, which was
447 halved if the loss did not improve over five iterations. The training
448 was continued until the loss did not decrease for ten consecutive
449 iterations (Figure 3), reducing the training time compared to relying
450 on the default reduction in the learning rate used by ADAM.



451

452 *Figure 3 The training process of the machine learning (ML) models. (a) The ML*
453 *model is trained on the five features from AMIGO with the recharge rate we want to*
454 *predict the groundwater response. (b) The ML models are based on the encoder*
455 *decoder architecture. Two variants of U-NET have skip connections from the*
456 *encoder directly to the decoder represented by the dashed line. The model weights*
457 *are iteratively updated during training using the ADAM optimiser to a minimise loss.*
458 *This loss is based on the mean squared error (MSE) between the (c) ML model*
459 *predictions and those from the numerical model AMIGO. The training iterations are*
460 *concluded when the loss does not reduce on an unseen validation set for ten*
461 *iterations. Basemap from OpenStreetMap-carto.*

## 2.3. Training the model

*Inputs to the ML models* - The relative significance of the phreatic aquifer's properties was evaluated using the permutation importance approach (Altmann et al., 2010). This method estimates the significance by evaluating the increase in error that occurs after permuting that property. This method was used to compare the importance of the five phreatic aquifer properties and eliminate irrelevant ones. The compared properties are: transmissivity, hydraulic resistance below the aquifer, DRN and RIV conductance, DRN and RIV stage relative to the groundwater head in the baseline scenario, and the ground height relative to the groundwater head in the baseline scenario. These properties are two dimensional and need to be summarised as tabular features before estimating their importance. Four tabular features are calculated from each 2D feature which are: (1) mean, (2) minimum, and (3) maximum values where the steady-state groundwater response was more than 1 cm and the (4) average value within a 50 m radius of the site. The two definitions of the area (where the response was more than 1 cm and 50 m from the site) were included to capture the influence of the geo-hydrological properties near the recharge site and away from the site. Three key characteristics of the response were used to assess the relevance of the features and to quantify the ML model performance: the area, the maximum, and the total groundwater response (Figure 4). The area of the response is defined as the area around the recharge site with more than 1 cm of groundwater response. The maximum response is the highest, and the total response is the volume of the aquifer saturated by a response of more than 1 cm. Based on the permutation importance, all five phreatic aquifer properties are used to train the model.

*Data preprocessing was performed to improve the representation of aquifers as inputs to the ML model. The 15 model layers in AMIGO are discontinuous and are often represented by thin, highly transmissive layers. For the input of the ML model, only the characteristics of the first aquifer were used. We defined the first aquifer by combining the layers until the resistance below it exceeds 200 days. This aquifer mostly consists of all 15 layers to the East and four layers towards the West. The resistance below the 15 layers was represented by the highest resistance between the layers in AMIGO (Figure 3A). The transmissivity of this aquifer is calculated based on the hydraulic conductivity and saturated thickness of the individual layers in the baseline scenario. The properties of the aquifers also exhibited right-skewed distributions with long tails, as evidenced by their interquartile ranges (Table 1). To improve the ML model's stability and performance, these properties were log-transformed and min-max scaled to 0 and 1. However, DRN and RIV stage and surface height relative to the average groundwater head were not transformed or scaled as they are linearly related to the maximum response, draining excess recharge.*

12

510 *Table 1 Range and interquartile range of input features from AMIGO, before scaling*
511 *and log-transforming some of the inputs.*

| Input | Min | First quartile | Median | Third quartile | Maximum | Scaled and log-transformed |
|---|---|---|---|---|---|---|
| Aquifer transmissivity (m²/day) | 0.2 | 920 | 1350 | 1785 | 5034 | ✓ |
| Aquitard resistance (day) | 200 | 4610 | 42650 | 171709 | 171709 | ✓ |
| DRN and RIV conductance (m²/day) | 0.002 | 7.6 | 10.0 | 18.1 | 5053.0 | ✓ |
| DRN and RIV stage relative to the baseline groundwater head (m) | -2.7 | -0.04 | 0.22 | 0.54 | 11.7 | |
| Surface level relative to the baseline groundwater head (m) | 0 | 0.9 | 1.2 | 1.7 | 49 | |

512 *Recharge scenarios* - The ML models are trained on the steady-state
513 groundwater response to additional aquifer recharge with a certain
514 rate applied for a certain site size, calculated by the numerical
515 model AMIGO. Scenarios with varying applied recharge rates, site
516 sizes, and locations were simulated to produce the data used to

517 train the ML models. The sites were selected using Latin Hypercube
518 Sampling (LHS) and Orthogonal Array-based Latin Hypercube
519 Sampling (OALHS) (Sándor and András, 2004). The recharge rates
520 applied to the topmost layer of the numerical model range from 5
521 mm/day to 25 mm/day, and the site sizes range from 0.01 km$^2$ to 1
522 km$^2$. Each site covers 16 to 1600 model cells (each model cell is
523 25x25 m). These ranges were selected to represent a complete
524 range of potential recharge sites. While there are no MAR projects
525 in the study area, there was a test site 8 km from the catchment. It
526 was 0.58 km$^2$ in size and recharged 5 mm/day during the growing
527 season (Tang et al., 2023). This site would fall within the range
528 considered. Internationally, recharge between 250 mm and 1500
529 mm is applied during the growing season, which equates to 1.4
530 mm/day to 8.3 mm/day (de Wit et al., 2022). While some sites
531 would fall below the range considered in this study, allowing for
532 higher recharge rates would enable identifying the maximum
533 potential recharge rate at the site. The recharge rate and site sizes
534 were selected using LHS to represent the entire range.

535 The effect of aquifer recharge is determined by the interplay of
536 multiple geohydrological properties that vary throughout the
537 catchment. While the geo-hydrological properties are the same for
538 all scenarios, we exposed the ML model to various combinations of
539 these properties by varying where the recharge is applied within the
540 model extent (Figure 2). The model extent covers 765 km$^2$, which
541 could consist of 75969 to 720 potential recharge sites. The location
542 of the sites was randomly selected to represent the entire model
543 extent in datasets of 100, 300, 500 or 1000 sites. Selecting the
544 location at random minimizes the potential for bias, ensuring better
545 model performance for all potential recharge sites. A similar
546 methodology is used to select locations in previous studies (He et
547 al., 2021; Taccari et al., 2022; Tao et al., 2022). Multiple sites were
548 simulated simultaneously while maintaining a minimum distance
549 between adjacent sites to reduce their interaction. Simulating
550 multiple sites limited the number of numerical model runs. We used
551 the OALHS method to ensure that samples are more evenly spaced,
552 even in multiple dimensions, unlike LHS. While OALHS ensures a
553 more uniform sampling, it does not guarantee a minimum distance
554 between adjacent points. To enforce this condition, adjacent points
555 are separated into groups, resulting in four numerical model
556 scenarios from each OALHS of x and y coordinates of the recharge
557 site's centres. Considering the dimensions of the model domain and
558 the minimum distance, 18 sites were sampled together and then
559 split into two groups of six and two groups of three sites. Multiple
560 OALHS were grouped to create datasets of various sizes that
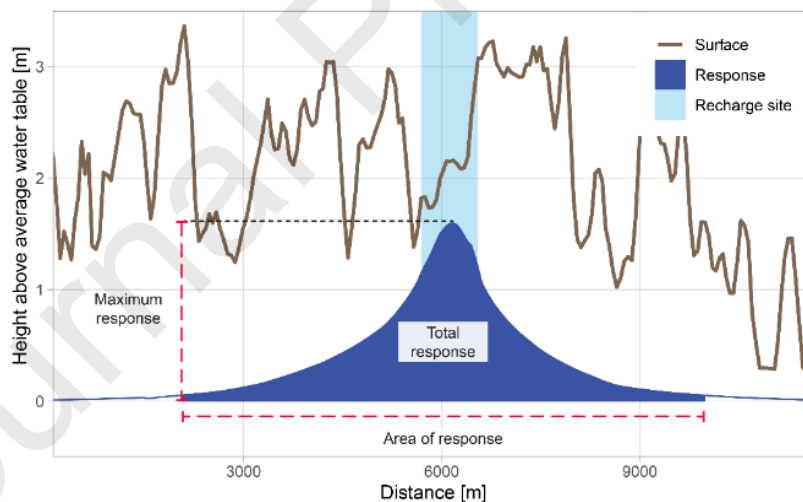561 represented the same sample distribution.

562 The results of the numerical model scenario runs were split into
563 three datasets: four training datasets, which were created through
564 resampling (with 1000, 500, 300, and 100 sites), a validation dataset
565 (100 sites), and a test dataset (200 sites). The OALHS samples were
566 maintained throughout the different datasets to ensure equal

567 representation. The recharge sites in each dataset are shown in
568 Appendix – A. The recharge sites in the training dataset with 1000
569 sites cover 364 km$^2$ representing 47.6% of the model extent. Of this,
570 47.3 km$^2$ overlaps with the test dataset. Although some sites in the
571 training dataset overlap with those in the testing dataset, the
572 recharge rate and the area of the sites differ between the sites.
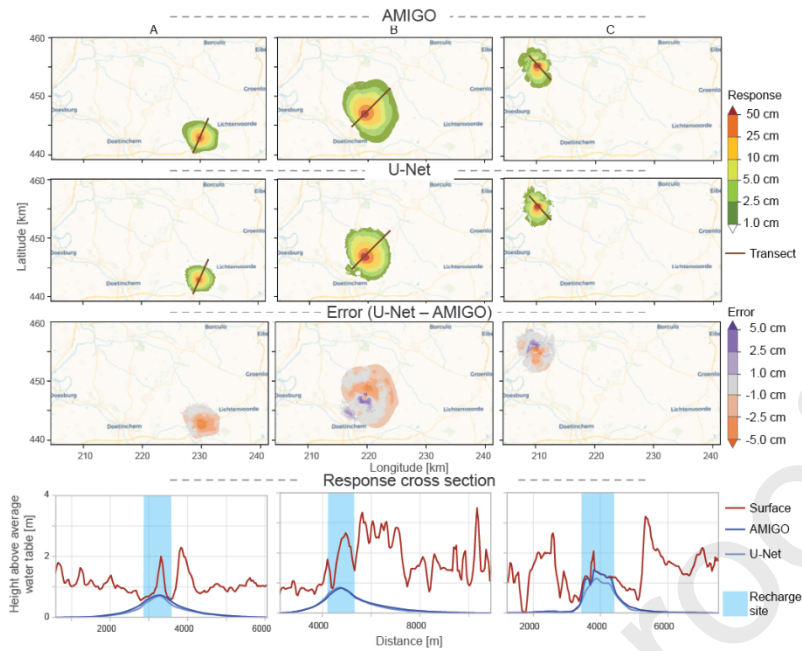
## 2.4. Analysis

574 The performance of the three ML models is assessed by comparing
575 their predictions of the three key characteristics, the maximum, the
576 area and the total response (Figure 4). The comparison uses the
577 Nash-Sutcliffe Efficiency (NSE) metric. The NSE measures the
578 model's ability to explain the variance in the observations and
579 ranges from -∞ to 1, with higher values implying a better predictive
580 ability.

581 While NSE describes the overall model's performance, it does not
582 account for systematic errors. The systemic errors across the range
583 of responses are represented in a scatter plot of the estimated key
584 characteristics from the best ML model and AMIGO. For this, we
585 considered scenarios with a constant recharge rate of 15mm/day
586 over recharge sites of 1 km$^2$ across the entire model domain.  The
587 key characteristics are also represented as maps that reveal the
588 interactions between the inputs and the resulting response.

589



590 *Figure 4 Cross-sectional view of a possible response of the groundwater to artificial*
591 *recharge. All heights are relative to the baseline groundwater head. The increase in*
592 *groundwater head (blue) is due to artificial recharge at the recharge site (light blue).*
593 *The brown line represents the surface elevation relative to the baseline*
594 *groundwater head. The maximum response, the area of the response and the total*
595 *response are the key characteristics used to quantify the model performance. The*
596 *vertical and horizontal axis are not symmetrical which exaggerates small changes in*
597 *groundwater depth and response.*

15

598

*Figure 5 Map view of the response estimated for three recharge sites (A, B, C),*
*represented in columns, by the numerical model AMIGO (top row) and by UNET*
*model (middle row) trained on 1000 input recharge sites. The difference between*
*the two is represented below as the error. The recharge sites are selected for their*
*asymmetric response caused by the interaction between the groundwater and the*
*surface water network (Groote Beek River and IJssel River). The bottom row*
*represents the cross sectional view of the response along the transects in the maps.*
*The vertical and horizontal axis in these cross sections are not symmetrical which*
*exaggerates small changes in groundwater depth and response. Basemap from*
*OpenStreetMap-carto.*

To showcase the advantages of the ML model, we undertook a
methodology aimed at determining the optimal location and
recharge rate for sites within the catchment area. This involved
simulating 7,722 recharge sites across the entire study area, each
covering an area of 10 hectares. The simulation included the
evaluation of eleven recharge rates ranging from 5 to 25 mm/day at
2 mm/day intervals for each site. Based on these simulations, we
created a database of 84942 responses among which the optimal
recharge sites can be identified. To identify these sites, we sought
locations exhibiting the highest response at a low recharge rate
based on the total volume of the response. Although related, this
target differs from the volume of water stored in the aquifer. To
estimate the extra volume of water which can be stored by the
artificial recharge, we multiplied the total volume of the response by
the specific yield of the phreatic aquifer. The specific yield used in
AMIGO for transient simulations is 0.15 *(Vreugdenhil, 2021)* which
corresponds to an aquifer composed of silt to medium sand
*(Johnson, 1967)*. This aquifer material type fits the description of the
Pleistocene sands in the catchment.

The assessment of these locations involved comparing their
response to a constant recharge rate of 25 mm/day. Subsequently,
the optimal recharge rate was discerned by identifying the minimum

16

631 recharge rate that achieved more than 80% of the maximum
632 response at each site. This comprehensive methodology allowed us
633 to systematically analyse and pinpoint the most effective locations
634 and recharge rates for artificial recharge within the catchment area
635 while demonstrating the benefits of the ML model.

## 3.  Results & Discussion

637 This section evaluates the performance of three ML models
638 (encoder-decoder, U-Net and attention U-net) in predicting the
639 steady-state groundwater head responses to artificial recharge,
640 generated by the numerical groundwater model, AMIGO. The best
641 performing ML model has captured the asymmetric responses to
642 the artificial recharge (Figure 5). This asymmetry is caused by the
643 interaction of the groundwater with the surface water network such
644 as, with rivers and drains. The surface water network drains part of
645 the groundwater response, hence limiting the response. Despite the
646 added complexity, the best ML model captured this interaction,
647 predicting the response outside the recharge site within ±10 cm. In
648 the following sections, the performance of the ML models is further
649 examined.

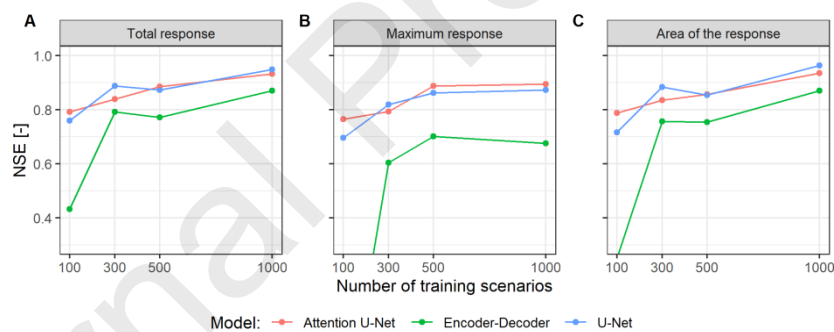### 3.1. Performance of the three machine learning models

652 All three ML models perform well when trained on 300 or more
653 recharge sites, indicated by high (NSE values (Figure 6). They
654 achieved a high NSE in comparing the total response and its area
655 despite the lower NSE for the maximum response.

656 Both U-Net and Attention U-Net models exhibited similar
657 performance and consistently outperformed the encoder-decoder
658 model. The variants of U-NET's outperformance could be due to the
659 increased number of model parameters and the significance of the
660 skip connections from the encoder to the decoder block in the U-
661 Net models (Figure 3). These skip connections allow the models to
662 capture spatially highly variable details in the input, such as DRN and
663 RIV properties. This conclusion is further supported by the encoder-
664 decoder model's worse performance at predicting the maximum
665 response, as high groundwater heads are strongly influenced by the
666 surface drainage network near them, which is better captured by
667 the U-Net models. This effect of the surface drainage network is
668 evident in Figure 5C, where the IJssel river (Figure 2) drains some of
669 the groundwater, causing an asymmetric response. Similarly, the
670 Grote Beek stream, southwest of the recharge site, causes a smaller
671 and steeper response (Figure 5B).

672 Attention U-Net learns to focus on important regions within the
673 input that help it predict the local response more accurately.
674 Contrary to its expected better accuracy, attention U-Net does not
675 have a significantly different NSE than U-Net. After accounting for
676 different training sizes, the true difference in NSE between the two

17

677  ML models is between -0.06 and 0.05 (95th percentile) based on
678  paired student's t-test. This result is counter-intuitive as the
679  response is highly localised and the ML models could gain from
680  focussing on selected parts of the input data. However, the
681  attention mechanism in attention U-Net's decoder block increases
682  the model's memory requirement, which we compensated for by
683  halving the number of filters in the convolution layers in the model.
684  Based on this, we can conclude that more filters greatly improve the
685  model performance, more than the advantages of the attention
686  layers. For models with a smaller extent, requiring less memory, it
687  could be more beneficial to train U-Net with more filters rather than
688  using Attention U-Net.

689  Furthermore, all three models improve with additional training data,
690  particularly for the area of the response and total response (Figure
691  6A and Figure 6C). Specifically, the U-Net model's NSE for the
692  predicted area increased from 0.71 to 0.96 with 1000 training sites
693  versus 100 sites, and the NSE value for the predicted total response
694  increased from 0.76 to 0.95 with additional training sites. However,
695  the NSE for the predicted maximum response did not consistently
696  improve with additional training data (Figure 6B). Additional training
697  sites improved the performance up to 500 sites, but the predicted
698  maximum response only marginally improved when doubling the
699  input to 1000 sites (NSE of U-Net from 0.86 to 0.87).
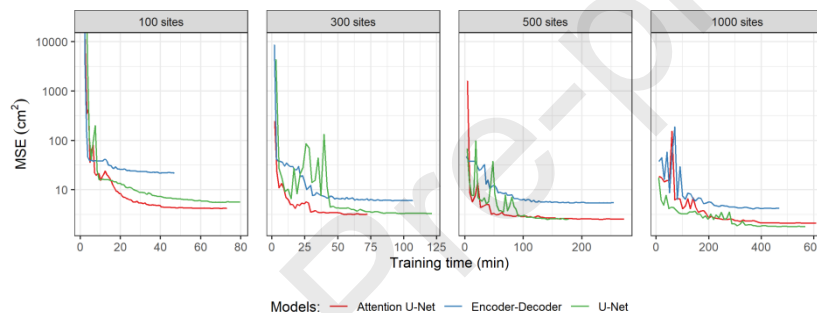


700

701  *Figure 6 Nash-Sutcliffe Efficiency (NSE) of the key characteristics from the*
702  *groundwater head response estimated by the machine learning models when*
703  *trained with an increasing number of training sites along the x-axis. A high NSE*
704  *(maximum of 1) indicated more accurate predictions. The three characteristics of*
705  *the response (total, maximum and area of the response) are represented in columns*
706  *(A, B and C).*

707  Another consideration when choosing the model is the training and
708  evaluation time. However, the training time is strongly dependent
709  on the initial values of the parameters in the ML model and hence
710  might not be perfectly reproduced. The initial parameters also
711  explain the initial error that improves during the training process
712  (Figure 7). The error does not steadily reduce during training and
713  often fluctuates, especially early into the training. This fluctuation is
714  likely due to a relatively high learning rate which was reduced when
715  the training stagnated. This learning rate reduced the overall
716  training time compared to relying on the ADAM optimiser's default
717  learning rate. The training process seems to be slowed by the

18

718  vanishing gradient problem, exacerbated by the sparse nature of the
719  response. The encoder-decoder model trained on 100 recharge sites
720  stagnated at this point and only predicted low responses. The model
721  needed more than 100 sites to train further.

722  The relative effect of the training size on the total training time
723  would likely be consistent in future training attempts. Additional
724  training sites linearly increase the training time, from 70 min when
725  trained on 100 sites to 10 hours for 1000 sites. Although it is a long
726  time, it is 'passive time' where no human interaction is required.
727  Each training iteration for the encoder-decoder model is shorter,
728  but it rarely outperformed the variants of U-Net (Figure 7). Between
729  the variants, Attention U-Net trained faster than U-Net for smaller
730  datasets with 100 sites and 300 sites and achieved lower validation
731  errors. This is likely due to the model's ability to learn regions to
732  focus on through training. However, U-Net can compensate for the
733  attention mechanism with additional training data and
734  outperformed Attention U-Net when trained on 1000 sites.



735

736  *Figure 7 Validation MSE that was tracked during the training process. The MSE is*
737  *calculated for an unseen set of recharge sites, validation set, different from the sites*
738  *used to train the model. Additional training sites improve the final model but also*
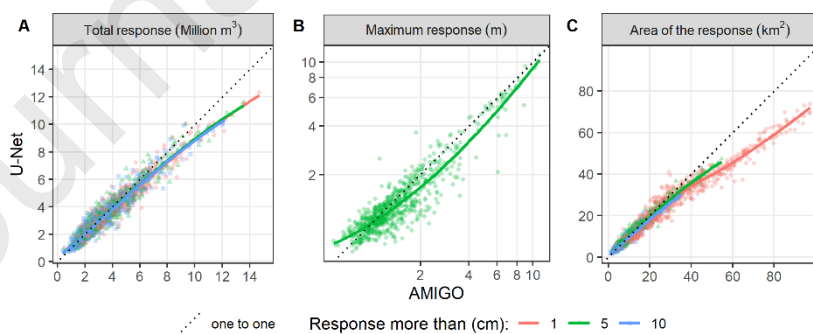739  *increase the training time.*

740  The evaluation time for the models ranges between 0.06 s to 0.43 s.
741  The average evaluation time for the three models ranged between
742  0.09 s and 0.11 s and varied significantly between the models
743  (Kruskal-Wallis test p-value < 0.01). However, this difference is not
744  of practical significance, especially when compared to the average
745  AMIGO run that took 1290 s (between 688 s and 2227 s). The
746  slowest ML model, U-Net, could evaluate 3000 scenarios during the
747  average time for a single scenario run in AMIGO.

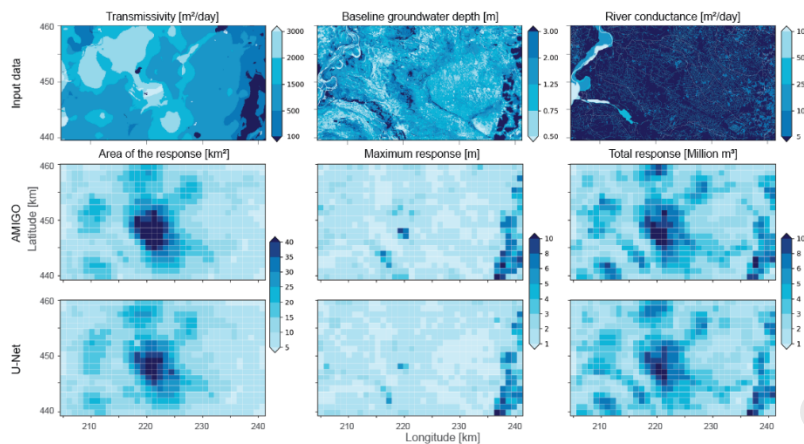## 3.2. Performance of the best model

749  The U-Net model trained on 1000 recharge sites is the best-
750  performing ML model with the highest NSE for predicting area and
751  total response. However, NSE does not account for systematic
752  errors. Figure 8 shows good agreement between the U-Net and
753  AMIGO estimates, but often U-Net underestimates the maximum
754  groundwater response.

755 Figure 5C is one such scenario where U-Net underestimates the
756 maximum response; limiting the response to the bottom of a local
757 depression at the recharge site. The recharge rate at the site
758 exceeds the maximum rate the groundwater can spread away from
759 the site, leading to the groundwater head reaching the surface and
760 seeping out through overland flow (cross-section of Figure 5C).
761 Although such a high recharge rate is not efficient at storing water in
762 the subsurface, the occurrence of overland flow would encourage a
763 redesign of the recharge site. AMIGO can capture this phenomenon,
764 but the response from U-Net does not reach the surface level. The
765 response from U-Net is limited by the deepest surface point,
766 resulting in a larger error at the recharge site (Figure 5C). However,
767 U-Net underestimates the response, which still suggests that the
768 site is inefficient at storing water and motivates redesigning the
769 recharge site. Additionally, this error has a minor impact on the
770 response away from the recharge site, where the response from
771 both AMIGO and U-Net are mostly within 7.2cm of each other (99th
772 percentile). This error is comparable to the responses in Figure 5A
773 and Figure 5B, 9.1cm and 5.4cm respectively.

774 The U-Net model shows a negative bias for high values in both the
775 total response and area of response. Specifically, the U-Net model
776 underestimates the response of the top ten sites with the highest
777 total response by 15% (see Figure 8A) and the area of the top ten
778 widest response by 13%. Notably, these results are only applicable
779 to responses more than 5cm. When including the smaller responses,
780 up to 1cm, the bias increases to 26% (Figure 8C). Interestingly,
781 increasing the lower limit to 10cm did not decrease the bias (13.1%
782 vs 13.0%), indicating that U-Net underestimates the small responses
783 and the bias increases for responses less than 5cm. Although the
784 total response is less sensitive to the minimum limit, it still increases
785 from 12% to 16% when considering responses less than 5cm.



786

787 *Figure 8 Scatter plot of the key characteristics of the response, estimated by U-Net*
788 *vs those from the numerical model, AMIGO. These results are for recharge sites*
789 *across the entire model domain, with 15mm/day recharge applied over 1 km². The*
790 *total response and area of the response were calculated for responses of more than*
791 *1cm, 5cm and 10cm to indicate the model's accuracy at predicting smaller*
792 *responses. The line is used to represent the trend in the scatter created from a local*
793 *polynomial regression fitting.*

794

*Figure 9 A comparison of the input data (top row) and the predictions of the three key characteristics (total, maximum and area of the response) from the numerical groundwater model (AMIGO, middle row) and our best machine learning model (U-Net trained on 1000 recharge sites, bottom row) for 1km² recharge sites with 15mm/day artificial recharge over the model domain.*

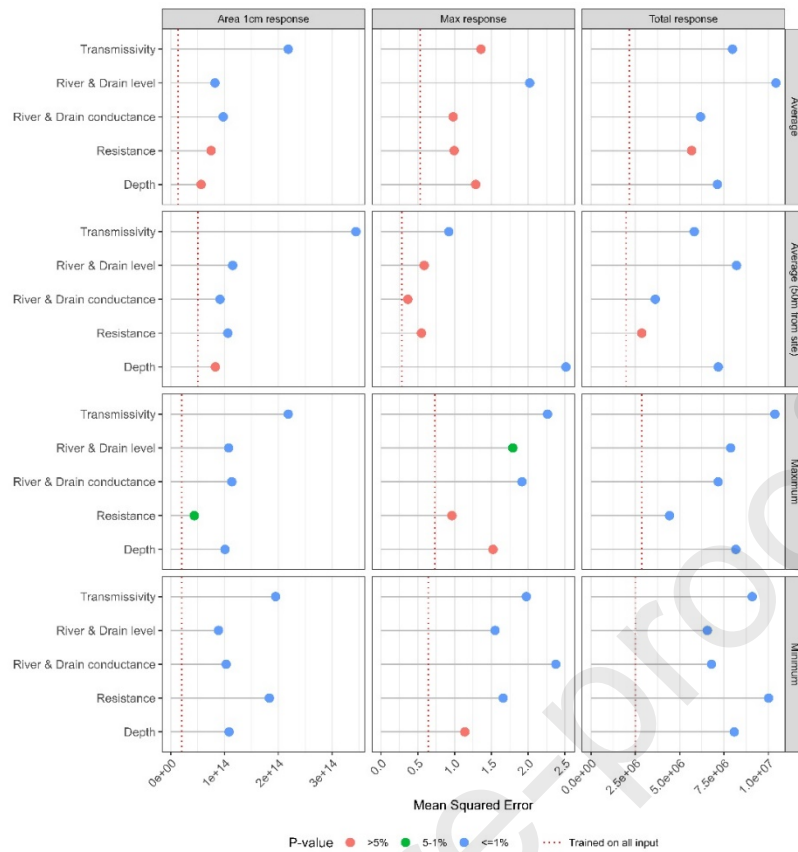### 3.3. Input features

The key response characteristics: area, maximum, and total response, depend on various hydro-geological inputs and their interaction. This interaction is evident in Figure 9, where the key characteristics are not directly related to any single hydro-geological input but a combination. U-Net could reproduce the spatial patterns of the key characteristics accurately, indicating that it has captured the effect of the interaction. Among the key characteristics, the maximum response has the most direct dependence on the groundwater depth below the ground surface and the DRN and RIV stage. These inputs limit the maximum response by draining some of the excess recharge. The recharge increases the groundwater head in the aquifer up to the drainage level. As the head increases above the drainage level, the groundwater is drained to the surface water network, depending on the head above the drain level and the drain conductance. This dependence is evident for sites at the elevated regions near the rivers. These rivers have a high conductance and hence more strongly limit the groundwater response. This dependency is also captured when estimating the importance of the inputs using the permutation importance approach (Figure 10). This approach suggests that the maximum response is significantly dependent on the average depth near the site, the maximum and minimum drain conductivity, transmissivity, and the minimum resistance.

The area of the response is the second key characteristic with a more direct relation to the input variables. The area depends on the aquifer's transmissivity, the minimum resistance between the aquifers, and the level and conductance of the surface drainage network (Figure 9) which is also reflected in the permutation importance (Figure 10). Higher transmissive aquifers allow for a faster flow of water away from the recharge site at a gentler

21

831  gradient. The faster flow and a gentler gradient result in the artificial
832  recharge providing water to a wider area. The dependency on the
833  surface drainage network can be explained by making a comparison
834  with groundwater abstraction. For groundwater abstractions, the
835  equation for leakage factor is related to the area around an
836  abstraction well where leakage occurs through the aquitard due to
837  the pumping in the aquifer below. Higher leakage factors indicate
838  that pumping would reduce the groundwater head in a wider area,
839  increasing the leakage in that area. Leakage factor ($\lambda$) is the square
840  root of the ratio of the aquifer's transmissivity (KD) and the aquitard
841  conductance (K'/D') above the aquifer: $\lambda = \sqrt{\dfrac{KD}{K'/D'}}$, where K and D
842  are the hydraulic conductivity and thickness of the layers. Phreatic
843  aquifers do not have an overlying aquitard; for these aquifers, the
844  properties of the surface drainage network are used instead (van
845  der Gaast et al., 2005). Besides the effect of the aquifer
846  transmissivity and resistance of the surface water network, the
847  permutation importance also suggests that the area of the response
848  depends on the minimum aquitard resistance below the aquifer
849  (Figure 10). However, the maximum and average resistance is only
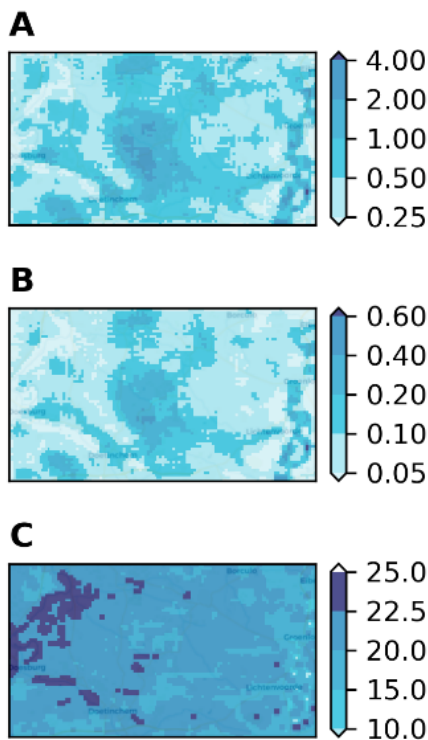850  significant up to a level of 5%.

851  The total response is the most complex and important key
852  characteristic of the response, related to the total volume of fresh
853  water stored using artificial recharge. It combines the other two key
854  characteristics, i.e. the maximum and the area of the response. The
855  total response also depends on the transmissivity and the surface
856  drainage network properties as they affect both the maximum and
857  the area of the response (Figure 9). Along with these inputs, the
858  total response depends on the average groundwater depth near the
859  site and the minimum aquitard resistance below the aquifer up to a
860  significance level of 1% (Figure 10). Based on these results, all five
861  features are necessary to ensure an adequate representation of the
862  system in the ML model.

863

*Figure 10 The permutation importance of the hydrological properties of the first*
*aquifer. This importance is the increase in the mean squared error at predicting*
*three performance indicators when the hydrological properties of the first aquifer*
*are randomized. The mean, minimum, and maximum values of the property where*
*the groundwater response was more than 1cm and the average of the property*
*within a 50m radius of the site were used to represent the hydrological properties*
*influencing the response at the site. The three performance indicators are (1) the*
*area of the groundwater response, (2) the maximum response, (3) total response. P-*
*values show the significance of the input characteristics in explaining the*
*performance indicators. The average, maximum or minimum of the hydrological*
*properties are important to explaining the key characteristics of the response up-to*

875 *a significance level of 0.01 and were hence included as inputs to the ML model.*



876

877 *Figure 11 Optimal recharge rate for 10 ha recharge sites across the entire study*
878 *area. The total volume of the response, in million m3, to recharge of 25mm/day*
879 *applied in sites of 10 ha is shown in A. However, this recharge rate is often*
880 *inefficient. B is the volume of water stored, in million m3, when recharging at a rate*
881 *that achieves at least 80% of the response at 25 mm/day. The corresponding*
882 *recharge rate, in mm/day, is in C.*

## 3.4. Applications

884 The ML model's efficiency, being 3000 times faster than the
885 numerical groundwater model, makes it suitable for various
886 applications requiring numerous steady-state model runs. For
887 instance, it can greatly benefit tasks like optimizing recharge rates,
888 determining the optimal size and location of recharge sites, and
889 comparing multiple locations rapidly. In cases where recharge
890 volume is predetermined, such as by regulatory mandates, the ML
891 model enables swift comparison of multiple locations. This
892 facilitates the evaluation of various combinations of recharge rates
893 and site areas, aiding in decision-making processes.

894 The bottom row of Figure 9 illustrates a notable example where the
895 key characteristics of 720 recharge sites were compared. The ML
896 model efficiently simulated these 720 recharge sites within 144
897 seconds, while the numerical model required 11 hours for the same
898 task. To enhance the speed of the numerical model runs, each run
899 simulated 6 equally spaced recharge sites, and five runs were
900 executed in parallel. This comparative analysis underscores the
901 substantial speed-ups achieved when using the ML model.

902    The results highlight specific regions within the catchment area,
903    particularly the center and eastern edges, as promising potential
904    recharge sites. Additionally, smaller regions near the northern and
905    southwestern edges of the model domain show promise. Figure 11A
906    depicts a similar analysis using the ML model, focusing on recharge
907    at a rate of 25 mm/day over 10 ha sites within the model domain.
908    The results reveal that, at this recharge rate, only the center and
909    eastern regions exhibit a high total response. This observation
910    suggests that different locations are more suitable at different
911    recharge rates.

912    To illustrate this point, we conducted a comprehensive comparison
913    involving the steady-state response of 7,722 recharge sites, each
914    covering an area of 10 hectares, across 11 recharge rates ranging
915    from 5 to 25 mm/day at 2 mm/day intervals. In total, the response
916    from 84,942 scenarios were predicted with the ML model in 980
917    seconds, which would have taken the numerical model 270 days
918    with the optimizations used to simulate 720 sites. This analysis
919    aimed to determine the minimum recharge rate that achieves 80%
920    of the highest total response for each site. Recharge sites located in
921    the eastern region of the catchment achieved a high total response
922    volume, saturating up to 4.35 million $m^3$ (Figure 11A), corresponding
923    to 0.65 million $m^3$ water stored (Figure 11B) at the optimal recharge
924    rate of 11 mm/day (Figure 11C). This effectiveness could be
925    attributed to the relatively low subsurface transmissivity, resulting
926    in a localized response to artificial recharge. Consequently, the
927    influence of streams and ditches away from the recharge site is
928    minimized. The high steady-state response achieved at a low
929    recharge rate makes this region emerge as a favourable location for
930    artificial recharge.

931    Conversely, the central portion of the model domain exhibits a
932    relatively high total response of 3 million $m^3$ while storing 0.45
933    million $m^3$ of water. This site benefits from a higher recharge rate of
934    23 mm/day (Figure 11C). Given the widespread response of these
935    recharge sites (Figure 9), they hold the potential to effectively raise
936    the groundwater level for the entire area, thereby enhancing water
937    availability for the broader natural environment. This underscores
938    the strategic importance of optimizing recharge rates based on the
939    specific characteristics of different regions to maximize the positive
940    impact on groundwater levels and ecosystem sustainability.

941    This analysis can readily incorporate variations in storage
942    coefficients across the model domain. By leveraging available data
943    on storage coefficients, we can optimize both stored water and
944    increases in groundwater head. While this integrated approach
945    would enable a comprehensive assessment of the groundwater
946    response and the alleviation of water stress on natural ecosystems
947    and the environment it is important to note that our current study
948    focuses on the steady-state response which is independent of the
949    storage coefficients and hence including the coefficient is beyond
950    the scope of this study. This focus allows us to delve deeply into the

25

951 system's long-term behaviour without the added complexity of
952 variable coefficients.

## 3.5. Steady-state vs transient scenarios

954 Steady-state scenarios depict the groundwater heads in a state of
955 equilibrium, where the inflows balance outflows without changes in
956 the storage within the cells. However, these scenarios assume no
957 changes in the boundary conditions throughout the simulation, such
958 as recharge, DRN, and RIV properties. These scenarios are thus not
959 intended to accurately reflect temporal dynamics, such as seasonal
960 variations in precipitation. Nevertheless, steady-state scenarios
961 provide valuable initial estimates, particularly for evaluating the
962 long-term effects of adaptation measures such as artificial recharge
963 when applying a constant recharge rate. Moreover, they require less
964 input data than transient scenarios and are faster to simulate than
965 transient scenarios. As a result, the data for training the ML model
966 are often available, making our technique applicable to more areas.
967 Having fewer input data that do not change during the simulation
968 facilitates precise attribution of the changes between scenarios to
969 specific inputs. This study leverages the benefits of steady-state
970 scenarios to demonstrate the applicability of the technique to
971 optimize artificial recharge sites.

972 Transient scenarios have the advantage that they can offer a more
973 detailed depiction, especially on the response of groundwater heads
974 and storage to artificial recharge in time, by accounting for the
975 dynamic nature of the system, based on which we can assess the
976 effect of seasonal variability on the system. Transient scenarios also
977 explicitly account for the changes in storage within each time step
978 due to the additional artificial recharge or due to seepage to the
979 surface water network. Understanding the effect of the geo-
980 hydrological properties that affect the changes in storage could
981 enhance the optimization of recharge site locations. Given the
982 successful development of an ML technique to mimic steady-state
983 conditions, as is done in the current study, the next step to develop
984 such an approach for transient conditions is warranted. It should be
985 noted however that successful implementation is not a given, as
986 complexity increases. This concerns e.g. differentiation between
987 periods of infiltration building up certain storage (autumn and
988 winter) and storage decay during summer seasons, for which
989 different ML approaches might be needed.

## 4. Conclusions

991 This study aims to understand the design choices for a machine
992 learning (ML) model to predict the steady-state groundwater
993 response to artificial recharge. It compares three state-of-the-art ML
994 models that best reproduce the response based on an identified
995 subset of the geo-hydrological data. The ML models were trained on
996 the results from a pre-calibrated numerical groundwater model to
997 reproduce the simulated response. In doing so, the response can be

998  estimated nearly instantly and help select appropriate artificial
999  recharge sites and optimise the sites. The ML model's performance
1000  was judged based on their performance at three key response
1001  characteristics: the maximum response, the area of the response
1002  and the total response.

1003  Three convolutional neural networks were trained, of which U-Net
1004  and Attention U-Net could accurately reproduce the response.
1005  These models contain skip connections that enable the model to
1006  capture spatially highly variable details in the inputs, such as DRN
1007  and RIV. Additionally, both these models have similar performance
1008  suggesting that the attention mechanism does not compensate for
1009  its memory requirement. With more available memory, training a U-
1010  NET with more filters could be more beneficial than opting for
1011  Attention U-NET. Both variants of U-NET achieved a high Nash
1012  Sutcliffe Efficiency (NSE) of 0.9 when trained on the results from 500
1013  recharge sites. Additional training sites improved the NSE to 0.96 at
1014  predicting the area of the response and the total response, while
1015  the maximum response did not show a marked improvement to
1016  additional data. However, additional data increases the computation
1017  time to generate the data and train the model, negating some of the
1018  benefits of the speed-up from the ML model. Despite the increased
1019  computation, the trained ML models could then be used to consider
1020  more scenarios, estimating the response within 0.24 s (95$^{th}$
1021  percentile), significantly faster than the numerical model, which
1022  took 1290 s. The slowest ML model, U-Net, could evaluate 3000
1023  scenarios during the average time for a single scenario run in
1024  AMIGO.

1025  Although the ML models trained in this study have a high NSE, they
1026  have their limitations. The models underestimate the maximum
1027  response in cases where groundwater levels reach the surface. Our
1028  best model is U-NET trained on 1000 sites; it limits the head to the
1029  deepest point at the recharge site (cross-section in Figure 5C). This
1030  error leads to underestimating the total response for scenarios with
1031  a high response. Despite this underestimation, the results do not
1032  impact the final recommendation that the scenario is sub-optimal
1033  and a similar response is possible with a lower recharge rate.
1034  Furthermore, the underestimation has a minor impact on the
1035  response away from the site or on the total response which is the
1036  most important characteristic to increase the water availability.
1037  Another limitation of the model is the lower accuracy in predicting
1038  small responses of less than 5cm. However, the smaller responses
1039  have a minor impact on the total response and hence should not
1040  affect the optimisation of the recharge sites.

1041  When training similar models, future work must decide between the
1042  geo-hydrological inputs that adequately represent the groundwater
1043  system. While the groundwater head response to phreatic aquifer
1044  recharge is mostly dependent on the properties of the phreatic
1045  aquifer itself, deeper aquifers do impact the response. The deeper
1046  aquifers have a diminishing impact on the flow which we addressed

27

1047 by combining the numerical model layers with a low resistance
1048 between them and focussing specifically on the first aquifer. Despite
1049 the potential for enhancing the model's accuracy by incorporating
1050 the properties of deeper aquifers, the ML models trained on the
1051 properties of the first aquifer could reproduce the steady-state
1052 response. Among the properties, we identified five crucial
1053 properties based on the results of the numerical model's scenarios
1054 and Altmann's permutation importance approach: transmissivity,
1055 resistance below the phreatic aquifer, depth to the groundwater,
1056 the water level in the surface water network and the network's
1057 hydraulic conductance to flow into the aquifer. Among these inputs,
1058 the transmissivity and surface water network properties are the
1059 most important as they impact all the key characteristics of the
1060 response. Considering the importance of these inputs, future
1061 research could focus on the effect of artificial recharge on these
1062 inputs. While the effect of higher transmissivity due to higher
1063 saturated thickness is incorporated in the numerical model
1064 simulations, the higher river stages due to greater flux to the river
1065 are not incorporated. A higher river stage would reduce the river
1066 flux which would increase the response. However, incorporating this
1067 would require generating the training data using a coupled surface
1068 water – groundwater model which is beyond the scope of this
1069 research.

1070 Fast models for specific tasks could prove an effective aid in
1071 designing good aquifer recharge sites. The speed-up could enable
1072 the water management authorities to consider many more
1073 scenarios in and around the selected catchment. The increased need
1074 for such an approach also follows from literature, e.g. from using
1075 ML-models to explain groundwater fluctuations (Sahoo et al., 2017)
1076 and the exploration of the influence of different uncertainties
1077 including future climate conditions while considering 1872 future
1078 scenarios (Miro et al., 2021). The approach could also motivate and
1079 justify the decisions to stakeholders improving support for water
1080 conservation. While this study does not demonstrate the model's
1081 performance in other regions, a similar model could best suit that
1082 region's challenges. The model in this study could serve as a starting
1083 point, and transfer learning techniques could be deployed, reducing
1084 the number of training scenarios needed and the training time.

1085 Finally, we identified challenges when covering a larger spatial
1086 extent by the model. The larger extent increases the spatial GPU
1087 memory required when training the machine learning model. The
1088 authors limited the size of the Attention U-Net to fit in the 16 G.B.
1089 available in NVIDIA Tesla T4 GPUs. Training an Attention U-Net with
1090 more filters could make it outperform U-Net. Similarly, the
1091 adversarial loss from generative adversarial networks (GANs) could
1092 further improve the model trained, but this required training an
1093 adversarial network alongside, increasing the memory overhead in
1094 the process.

1095 The models in this study focus on the groundwater response within
1096 the Baakse Beek catchment in the Netherlands. Future researchers
1097 could focus on training a single model for different locations, in
1098 order to investigate to what extent an ML model could be generally
1099 applicable and usable in catchments with sparse data. However, a
1100 similar extent must be maintained to ensure it can predict the entire
1101 spatial extent of the response. Furthermore, the current model is
1102 limited to steady-state scenarios, and considering the response's
1103 evolution during dryer periods could influence design choices.
1104 Groundwater heads are deeper during dryer periods, increasing the
1105 potential response to MAR. The groundwater fluctuations near the
1106 recharge site are sensitive to the storage coefficient of the
1107 surrounding aquifer which is not considered in steady-state
1108 groundwater response. The U-Net trained in this study may be
1109 extended for more complex scenarios and can be used to capture
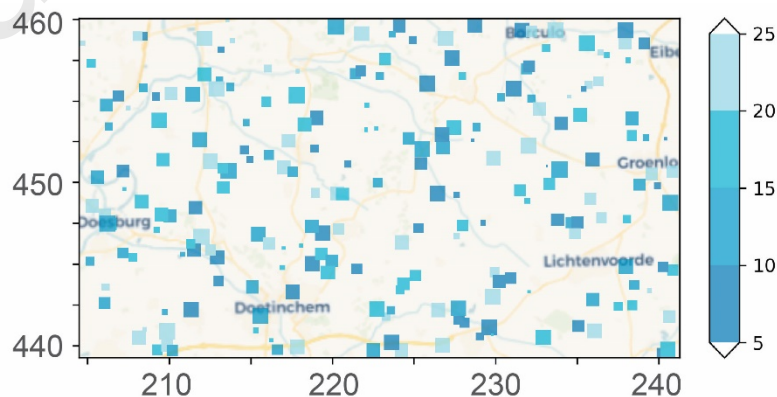1110 the effect of other geo-hydrological properties.

## 1111 5. Acknowledgements

## 1117 6. Declaration of generative AI and AI-assisted
## 1118 technologies in the writing process

1119 During the preparation of this work the authors used ChatGPT-3 for
1120 the sole purpose of improving the clarity of the text within this
1121 work. After using this tool, the authors reviewed and edited the
1122 content as needed and take full responsibility for the content of the
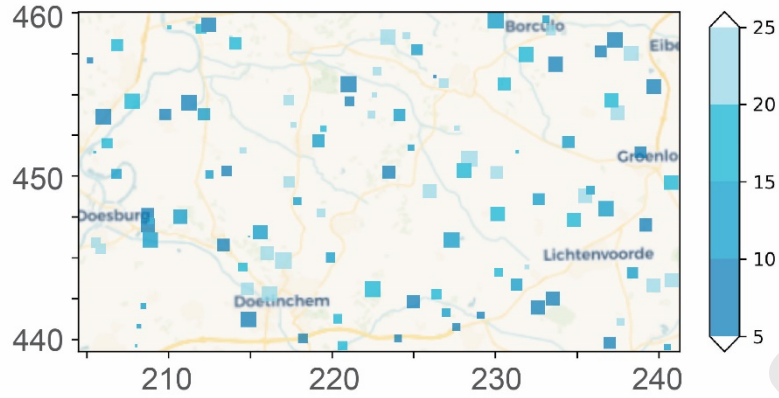1123 publication.

## 1124 Appendix – A: Site locations and recharge rates for all
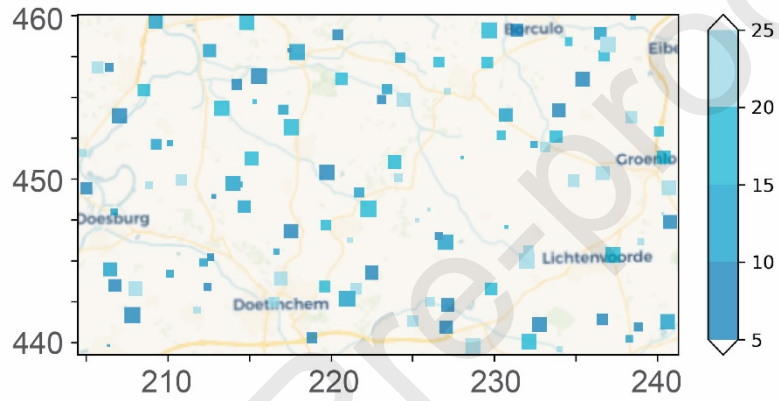## 1125 sites in the datasets



1126

1127 *Figure A2 Recharge rate at all the recharge sites in the Testing dataset. The three*
1128 *ML models are compared on their performance at predicitng the response to the*
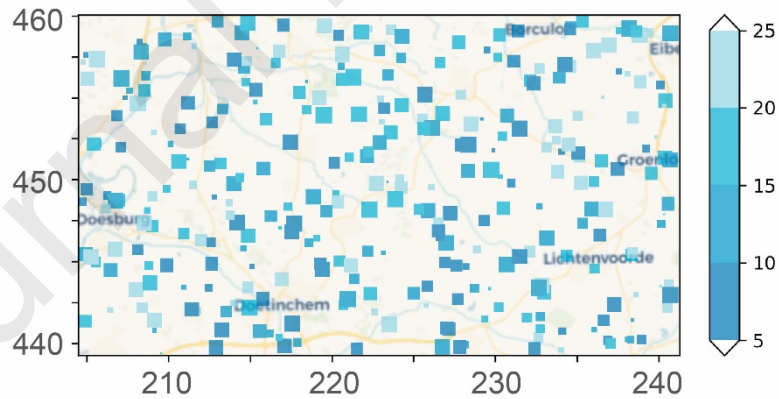
1129 *recharge sites in this dataset.*



1130

1131 *Figure A3 Recharge rate at all the recharge sites in the Validation dataset. This*
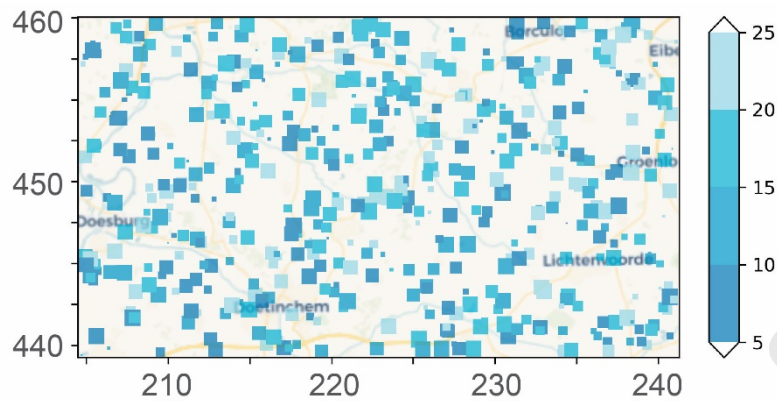1132 *dataset is used to track model performance during training.*



1133

1134 *Figure A4 Recharge rate at all the recharge sites in the Training dataset with 100*
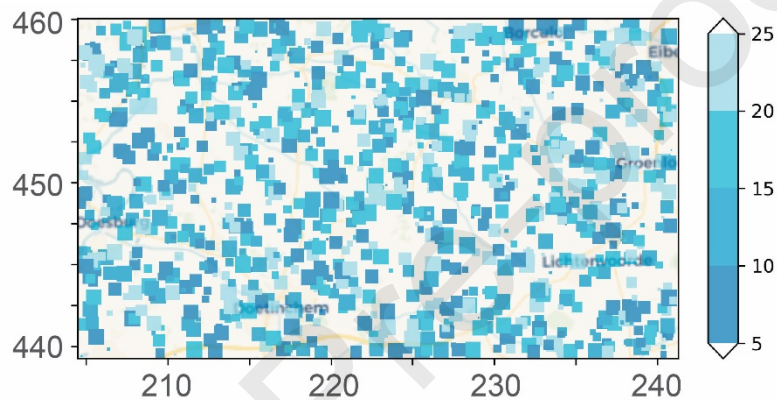1135 *sites*



1136

30

1137 *Figure A5 Recharge rate at all the recharge sites in the Training dataset with 300*
1138 *sites*



1139

1140 *Figure A6 Recharge rate at all the recharge sites in the Train dataset with 500 sites*



1141

1142 *Figure A7 Recharge rate at all the recharge sites in the Train dataset with 1000 sites*

1143 ## 7. References

1144 Aalbers, E.E., van Meijgaard, E., Lenderink, G., de Vries, H., van den
1145          Hurk, B.J.J.M., 2023. The 2018 west-central European
1146          drought projected in a warmer climate: how much drier can
1147          it get? Nat. Hazards Earth Syst. Sci. 23, 1921–1946.
1148          https://doi.org/10.5194/nhess-23-1921-2023

1149 Ahmadalipour, A., Moradkhani, H., Castelletti, A., Magliocca, N.,
1150          2019. Future drought risk in Africa: Integrating vulnerability,
1151          climate change, and population growth. Sci. Total Environ.
1152          662, 672–686.
1153          https://doi.org/10.1016/j.scitotenv.2019.01.278

1154 Altmann, A., Toloşi, L., Sander, O., Lengauer, T., 2010. Permutation
1155          importance: a corrected feature importance measure.
1156          Bioinformatics 26, 1340–1347.
1157          https://doi.org/10.1093/bioinformatics/btq134

1158 Asher, M.J., Croke, B.F.W., Jakeman, A.J., Peeters, L.J.M., 2015. A
1159          review of surrogate models and their application to
1160          groundwater modeling: SURROGATES OF GROUNDWATER

MODELS. Water Resour. Res. 51, 5957–5973.
https://doi.org/10.1002/2015WR016967

Balting, D.F., AghaKouchak, A., Lohmann, G., Ionita, M., 2021.
    Northern Hemisphere drought risk in a warming climate. Npj
    Clim. Atmospheric Sci. 4, 1–13.
    https://doi.org/10.1038/s41612-021-00218-2

Bartholomeus, R.P., Wiel, K. van der, Loon, A.F. van, Huijgevoort,
    M.H.J. van, Vliet, M.T.H. van, Mens, M., Geffen, S.M.,
    Wanders, N., Pot, W., 2023. Managing water across the
    flood–drought spectrum: Experiences from and challenges
    for the Netherlands. Camb. Prisms Water 1, e2.
    https://doi.org/10.1017/wat.2023.4

Bishop, C.M., 2006. Pattern recognition and machine learning,
    Information science and statistics. Springer, New York.

Boyce, S.E., Nishikawa, T., Yeh, W.W.-G., 2015. Reduced order
    modeling of the Newton formulation of MODFLOW to solve
    unconfined groundwater flow. Adv. Water Resour. 83, 250–
    262. https://doi.org/10.1016/j.advwatres.2015.06.005

Brakkee, E., Van Huijgevoort, M.H.J., Bartholomeus, R.P., 2022.
    Improved understanding of regional groundwater drought
    development through time series modelling: the 2018–2019
    drought in the Netherlands. Hydrol. Earth Syst. Sci. 26, 551–
    569. https://doi.org/10.5194/hess-26-551-2022

Brunton, S.L., Noack, B.R., Koumoutsakos, P., 2020. Machine
    Learning for Fluid Mechanics. Annu. Rev. Fluid Mech. 52,
    477–508. https://doi.org/10.1146/annurev-fluid-010719-
    060214

Casanova, J., Devau, N., Pettenati, M., 2016. Managed Aquifer
    Recharge: An Overview of Issues and Options, in: Jakeman,
    A.J., Barreteau, O., Hunt, R.J., Rinaudo, J.-D., Ross, A. (Eds.),
    Integrated Groundwater Management: Concepts,
    Approaches and Challenges. Springer International
    Publishing, Cham, pp. 413–434.
    https://doi.org/10.1007/978-3-319-23576-9_16

Castle, S.L., Thomas, B.F., Reager, J.T., Rodell, M., Swenson, S.C.,
    Famiglietti, J.S., 2014. Groundwater depletion during
    drought threatens future water security of the Colorado
    River Basin. Geophys. Res. Lett. 41, 5904–5911.
    https://doi.org/10.1002/2014GL061055

de Wit, J.A., van Dam, J.C., Ritsema, C.J., Bartholomeus, R.P., van
    den Eertwegh, G., 2022. Ontwikkeling van
    drainagesystemen: Water afvoeren - vasthouden -
    aanvullen. Stromingen Vakbl. Voor Hydrol. 28, 45–56.

32

1204    de Wit, J., Ritsema, C.K., van Dam, J.C., Van Den Eertwegh, G.A.P.H.,
1205            Bartholomeus, R.P., 2022. Development of subsurface
1206            drainage systems: Discharge – retention – recharge. Agric.
1207            Water Manag. 269, 107677.
1208            https://doi.org/10.1016/j.agwat.2022.107677

1209    Dey, S., Dhar, A., 2020. On proper orthogonal decomposition (POD)
1210            based reduced-order modeling of groundwater flow through
1211            heterogeneous porous media with point source singularity.
1212            Adv. Water Resour. 144, 103703.
1213            https://doi.org/10.1016/j.advwatres.2020.103703

1214    Dillon, P., Fernández Escalante, E., Megdal, S.B., Massmann, G.,
1215            2020. Managed Aquifer Recharge for Water Resilience.
1216            Water 12, 1846. https://doi.org/10.3390/w12071846

1217    Dillon, P., Stuyfzand, P., Grischek, T., Lluria, M., Pyne, R.D.G., Jain,
1218            R.C., Bear, J., Schwarz, J., Wang, W., Fernandez, E., Stefan,
1219            C., Pettenati, M., van der Gun, J., Sprenger, C., Massmann,
1220            G., Scanlon, B.R., Xanke, J., Jokela, P., Zheng, Y., Rossetto, R.,
1221            Shamrukh, M., Pavelic, P., Murray, E., Ross, A., Bonilla
1222            Valverde, J.P., Palma Nava, A., Ansems, N., Posavec, K., Ha,
1223            K., Martin, R., Sapiano, M., 2019. Sixty years of global
1224            progress in managed aquifer recharge. Hydrogeol. J. 27, 1–
1225            30. https://doi.org/10.1007/s10040-018-1841-z

1226    Harbaugh, A.W., 2005. MODFLOW-2005, The U.S. Geological Survey
1227            Modular Ground-Water Model—the Ground-Water Flow
1228            Process (Techniques and Methods), Techniques and
1229            Methods.

1230    Hartog, N., Stuyfzand, P., 2017. Water Quality Considerations on the
1231            Rise as the Use of Managed Aquifer Recharge Systems
1232            Widens. Water 9, 808. https://doi.org/10.3390/w9100808

1233    He, T., Wang, N., Zhang, D., 2021. Theory-guided full convolutional
1234            neural network: An efficient surrogate model for inverse
1235            problems in subsurface contaminant transport. Adv. Water
1236            Resour. 157, 104051.
1237            https://doi.org/10.1016/j.advwatres.2021.104051

1238    Hijma, M., 2017. Geology of the Dutch coast (No. 1220040- 007-
1239            ZKS- 0003).

1240    Ioffe, S., Szegedy, C., 2015. Batch Normalization: Accelerating Deep
1241            Network Training by Reducing Internal Covariate Shift.

1242    Johnson, A.I., 1967. Specific yield: compilation of specific yields for
1243            various materials (Report No. 1662D), Water Supply Paper.
1244            Washington, D.C. https://doi.org/10.3133/wsp1662D

1245    Kim, S., Melby, J.A., Nadal-Caraballo, N.C., Ratcliff, J.J., 2015. A time-
1246            dependent surrogate model for storm surge prediction

33

1247 based on an artificial neural network using high-fidelity
1248 synthetic hurricane modeling. Nat. Hazards 76, 565–585.

1249 Kingma, D.P., Ba, J., 2014. Adam: A Method for Stochastic
1250 Optimization. https://doi.org/10.48550/ARXIV.1412.6980

1251 Kutz, J.N., Brunton, S.L., 2022. Parsimony as the ultimate regularizer
1252 for physics-informed machine learning. Nonlinear Dyn. 107,
1253 1801–1817. https://doi.org/10.1007/s11071-021-07118-3

1254 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521,
1255 436–444. https://doi.org/10.1038/nature14539

1256 Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-Based
1257 Learning Applied to Document Recognition. Proc. IEEE 86,
1258 47.

1259 Lehner, F., Coats, S., Stocker, T.F., Pendergrass, A.G., Sanderson,
1260 B.M., Raible, C.C., Smerdon, J.E., 2017. Projected drought
1261 risk in 1.5°C and 2°C warmer climates. Geophys. Res. Lett.
1262 44, 7419–7428. https://doi.org/10.1002/2017GL074117

1263 Lerman, S., Venuto, C., Kautz, H., Xu, C., 2021. Explaining Local,
1264 Global, And Higher-Order Interactions In Deep Learning, in:
1265 2021 IEEE/CVF International Conference on Computer Vision
1266 (ICCV). Presented at the 2021 IEEE/CVF International
1267 Conference on Computer Vision (ICCV), IEEE, Montreal, QC,
1268 Canada, pp. 1204–1213.
1269 https://doi.org/10.1109/ICCV48922.2021.00126

1270 Lu, L., Shin, Y., Su, Y., Karniadakis, G.E., 2020. Dying ReLU and
1271 Initialization: Theory and Numerical Examples. Commun.
1272 Comput. Phys. 28, 1671–1706.
1273 https://doi.org/10.4208/cicp.OA-2020-0165

1274 Maas, A.L., Hannun, A.Y., Ng, A.Y., 2013. Rectifier Nonlinearities
1275 Improve Neural Network Acoustic Models. Presented at the
1276 International Conference on Machine Learning, JMLR:
1277 W&CP, Atlanta, Georgia, USA.

1278 Malik, A., Bhagwat, A., 2021. Modelling groundwater level
1279 fluctuations in urban areas using artificial neural network.
1280 Groundw. Sustain. Dev. 12, 100484.
1281 https://doi.org/10.1016/j.gsd.2020.100484

1282 Miro, M.E., Groves, D., Tincher, B., Syme, J., Tanverakul, S., Catt, D.,
1283 2021. Adaptive water management in the face of
1284 uncertainty: Integrating machine learning, groundwater
1285 modeling and robust decision making. Clim. Risk Manag. 34,
1286 100383. https://doi.org/10.1016/j.crm.2021.100383

1287 Mo, S., Zabaras, N., Shi, X., Wu, J., 2019. Deep Autoregressive Neural
1288 Networks for High-Dimensional Inverse Problems in

34

1289          Groundwater Contaminant Source Identification. Water
1290          Resour. Res. 55, 3856–3881.
1291          https://doi.org/10.1029/2018WR024638

1292 Müller, J., Park, J., Sahu, R., Varadharajan, C., Arora, B., Faybishenko,
1293          B., Agarwal, D., 2021. Surrogate optimization of deep neural
1294          networks for groundwater predictions. J. Glob. Optim. 81,
1295          203–231. https://doi.org/10.1007/s10898-020-00912-0

1296 Newman, A., 1996. Model Reduction via the Karhunen-Loeve
1297          Expansion Part I: An Exposition (Technical Research Report
1298          No. T.R. 96-32). Institute for Systems Research, University of
1299          Maryland.

1300 Oktay, O., Schlemper, J., Folgoc, L.L., Lee, M., Heinrich, M., Misawa,
1301          K., Mori, K., McDonagh, S., Hammerla, N.Y., Kainz, B.,
1302          Glocker, B., Rueckert, D., 2018. Attention U-Net: Learning
1303          Where to Look for the Pancreas.

1304 Papadopoulos, V., Soimiris, G., Giovanis, D.G., Papadrakakis, M.,
1305          2018. A neural network-based surrogate model for carbon
1306          nanotubes with geometric nonlinearities. Comput. Methods
1307          Appl. Mech. Eng. 328, 411–430.
1308          https://doi.org/10.1016/j.cma.2017.09.010

1309 Philip, S.Y., Kew, S.F., Van Der Wiel, K., Wanders, N., Jan Van
1310          Oldenborgh, G., 2020. Regional differentiation in climate
1311          change induced drought trends in the Netherlands. Environ.
1312          Res. Lett. 15, 094081. https://doi.org/10.1088/1748-
1313          9326/ab97ca

1314 Pronk, G.J., Stofberg, S.F., Van Dooren, T.C.G.W., Dingemans,
1315          M.M.L., Frijns, J., Koeman-Stein, N.E., Smeets, P.W.M.H.,
1316          Bartholomeus, R.P., 2021. Increasing Water System
1317          Robustness in the Netherlands: Potential of Cross-Sectoral
1318          Water Reuse. Water Resour. Manag. 35, 3721–3735.
1319          https://doi.org/10.1007/s11269-021-02912-5

1320 Rakovec, O., Samaniego, L., Hari, V., Markonis, Y., Moravec, V.,
1321          Thober, S., Hanel, M., Kumar, R., 2022. The 2018–2020
1322          Multi-Year Drought Sets a New Benchmark in Europe. Earths
1323          Future 10, e2021EF002394.
1324          https://doi.org/10.1029/2021EF002394

1325 Sahoo, S., Russo, T.A., Elliott, J., Foster, I., 2017. Machine learning
1326          algorithms for modeling groundwater level changes in
1327          agricultural regions of the U.S. Water Resour. Res. 53, 3878–
1328          3895. https://doi.org/10.1002/2016WR019933

1329 Sándor, Z., András, P., 2004. Alternative sampling methods for
1330          estimating multivariate normal probabilities. J. Econom.

1331  120, 207–234. https://doi.org/10.1016/S0304-
1332  4076(03)00212-4

1333  Sevink, J., Koopman, S., 2020. Maximum Holocene groundwater
1334  levels and associated extension of peat in the border zone of
1335  'Het Gooi' (the Netherlands): a reconstruction based on the
1336  study of soil transects. Neth. J. Geosci. 99, e7.
1337  https://doi.org/10.1017/njg.2020.7

1338  Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I.,
1339  Salakhutdinov, R., 2014. Dropout: A Simple Way to Prevent
1340  Neural Networks from Overfitting. J. Mach. Learn. Res. 15,
1341  1929–1958.

1342  Stanko, Z.P., Boyce, S.E., Yeh, W.W.-G., 2016. Nonlinear model
1343  reduction of unconfined groundwater flow using POD and
1344  DEIM. Adv. Water Resour. 97, 130–143.
1345  https://doi.org/10.1016/j.advwatres.2016.09.005

1346  Taccari, M.L., Nuttall, J., Chen, X., Wang, H., Minnema, B., Jimack,
1347  P.K., 2022. Attention U-Net as a surrogate model for
1348  groundwater prediction. Adv. Water Resour. 163.
1349  https://doi.org/10.1016/j.advwatres.2022.104169

1350  Tang, D.W.S., Van der Zee, S.E.A.T.M., Narain-Ford, D.M., van den
1351  Eertwegh, G.A.P.H., Bartholomeus, R.P., 2023. Managed
1352  phreatic zone recharge for irrigation and wastewater
1353  treatment. J. Hydrol. 626, 130208.
1354  https://doi.org/10.1016/j.jhydrol.2023.130208

1355  Tao, H., Hameed, M.M., Marhoon, H.A., Zounemat-Kermani, M.,
1356  Heddam, S., Kim, S., Sulaiman, S.O., Tan, M.L., Sa'adi, Z.,
1357  Mehr, A.D., Allawi, M.F., Abba, S.I., Zain, J.M., Falah, M.W.,
1358  Jamei, M., Bokde, N.D., Bayatvarkeshi, M., Al-Mukhtar, M.,
1359  Bhagat, S.K., Tiyasha, T., Khedher, K.M., Al-Ansari, N.,
1360  Shahid, S., Yaseen, Z.M., 2022. Groundwater level prediction
1361  using machine learning models: A comprehensive review.
1362  Neurocomputing 489, 271–308.
1363  https://doi.org/10.1016/j.neucom.2022.03.014

1364  Thatch, L.M., Gilbert, J.M., Maxwell, R.M., 2020. Integrated
1365  Hydrologic Modeling to Untangle the Impacts of Water
1366  Management During Drought. Groundwater 58, 377–391.
1367  https://doi.org/10.1111/gwat.12995

1368  Thomas, B.F., Famiglietti, J.S., 2019. Identifying Climate-Induced
1369  Groundwater Depletion in GRACE Observations. Sci. Rep. 9,
1370  4124. https://doi.org/10.1038/s41598-019-40155-y

1371  van den Eertwegh, G., Bartholomeus, R., de Louw, P., Witte, F., van
1372  Dam, J., van Deijl, D., Hoefsloot, P., van Huijgevoort, M.,
1373  Hunink, J., America, I., Pouwels, J., de Wit, J., 2020. Droogte

1374          in zandgebieden van Zuid-, Midden- en Oost-Nederland: Het
1375          verhaal: analyse van droogte 2018 en 2019 en tussentijdse
1376          bevindingen. KWR.

1377 van der Gaast, J.W.J., Massop, H.Th.L., Heuvelink, G.B.M., 2005.
1378          Monitpring van verdroging; Methodische aspected van
1379          meetnetoptimalisatie (No. 1102). Alterra, Wageningen.

1380 van der Wiel, K., Lenderink, G., de Vries, H., 2021. Physical storylines
1381          of future European drought events like 2018 based on
1382          ensemble climate modelling. Weather Clim. Extrem. 33,
1383          100350. https://doi.org/10.1016/j.wace.2021.100350

1384 Vermeulen, P.T.M., Heemink, A.W., Te Stroet, C.B.M., 2004.
1385          Reduced models for linear groundwater flow models using
1386          empirical orthogonal functions. Adv. Water Resour. 27, 57–
1387          69. https://doi.org/10.1016/j.advwatres.2003.09.008

1388 Vermeulen, P.T.M., Minnema, B., Roelofsen, F.J., 2021. iMOD User
1389          Manual version 5.3. Deltares manual.

1390 Vreugdenhil, I., 2021. Modelverbetering AMIGO 3.1 (No.
1391          D10043782:14). Arcadis Nederland B.V.

1392 Wang, C., Duan, Q., Gong, W., Ye, A., Di, Z., Miao, C., 2014. An
1393          evaluation of adaptive surrogate modeling based
1394          optimization with two benchmark problems. Environ.
1395          Model. Softw. 60, 167–179.
1396          https://doi.org/10.1016/j.envsoft.2014.05.026

1397 Weber, T., Corotan, A., Hutchinson, B., Kravitz, B., Link, R., 2019.
1398          Technical Note: Deep Learning for Creating Surrogate
1399          Models of Precipitation in Earth System Models (preprint).
1400          Clouds and Precipitation/Atmospheric
1401          Modelling/Troposphere/Physics (physical properties and
1402          processes). https://doi.org/10.5194/acp-2019-85

1403 Witte, J.P.M., Runhaar, J., Bartholomeus, R.P., Fujita, Y., Hoefsloot,
1404          P., Kros, J., Mol, J., de Vries, W., 2018. De waterwijzer
1405          natuur: instrumentarium voor kwantificeren van effecten
1406          van waterbeheer en klimaat op terrestrische natuur, Stowa
1407          rapport. Stowa.

1408 **Highlights**

1409 • U-Net accurately reproduces the groundwater response to artificial
1410 recharge
1411 • Inputs include properties of the first aquifer, drainage network, and
1412 recharge rate
1413 • Transmissivity and surface water networks significantly impact the
1414 response

1415• U-Net representing a 3000-fold speed up compared to gridded
1416   groundwater model
1417• Minor benefits from more than 500 recharge sites for training

1418