This is the accepted manuscript version of the contribution published as:

Cai, H., Liu, S., Shi, H., Zhou, Z., **Jiang, S.**, Babovic, V. (2022): Toward improved lumped groundwater level predictions at catchment scale: Mutual integration of water balance mechanism and deep learning method *J. Hydrol.* **613, Part B**, art. 128495

The publisher's version is available at:

http://dx.doi.org/10.1016/j.jhydrol.2022.128495

1	Toward improved lumped groundwater level predictions at catchment
2	scale: mutual integration of water balance mechanism and deep learning
3	method
4	
5	Hejiang Cai ^{1,2,3} , Suning Liu ⁴ , Haiyun Shi ^{2,3,*} , Zhaoqiang Zhou ^{2,3} , Shijie Jiang ⁵ , Vladan
6	Babovic ¹
7	
8	¹ Department of Civil and Environmental Engineering, National University of Singapore,
9	Singapore
10	² State Environmental Protection Key Laboratory of Integrated Surface Water-Groundwater
11	Pollution Control, School of Environmental Science and Engineering, Southern University of
12	Science and Technology, Shenzhen, Guangdong, China
13	³ Guangdong Provincial Key Laboratory of Soil and Groundwater Pollution Control, School of
14	Environmental Science and Engineering, Southern University of Science and Technology,
15	Shenzhen, Guangdong, China
16	⁴ Center for Climate Physics, Institute for Basic Science, Busan, Republic of Korea
17	⁵ Department of Computational Hydrosystems, Helmholtz Centre for Environmental Research,
18	Leipzig, Germany
19	
20	Corresponding author: Dr. Haiyun Shi (shihy@sustech.edu.cn)
21	

22 Abstract

Model development in groundwater simulation and physics informed deep learning (DL) has been 23 advancing separately with limited integration. This study develops a general hybrid model for 24 groundwater level (GWL) simulations, wherein water balance-based groundwater processes are 25 embedded as physics constrained recurrent neural layers into prevalent DL architectures. Because 26 of the automatic parameterizing process, physics-informed deep learning algorithm (DLA) equips 27 the hybrid model with enhanced abilities of inferring geological structures of catchment and 28 unobserved groundwater-related processes implicitly. The main purposes of this study are: 1) to 29 30 explore an optimized data-driven method as alternative to complicated groundwater models; 2) to improve the awareness of hydrological knowledge of DL model for lumped GWL simulation; and 31 3) to explore the lumped data-driven groundwater models for cross-region applications. The 91 32 illustrative cases of GWL modeling across the middle eastern continental United States (CONUS) 33 demonstrate that the hybrid model outperforms the pure DL models in terms of prediction 34 accuracy, generality, and robustness. More specifically, the hybrid model outperforms the pure DL 35 models in 78% of catchments with the improved $\Delta NSE=0.129$. Meanwhile, the hybrid model 36 simulates more stably with different input strategies. This study reveals the superiority and 37 38 powerful simulation ability of the DL model with physical constraints, which increases trust in data-driven approaches on groundwater modellings. 39

40 Key words

41 Groundwater level predictions; Water balance mechanism; Deep learning; Catchment scale

1

Highlights

2 •	A novel hybrid	model for simulating	groundwater level	was developed
-----	----------------	----------------------	-------------------	---------------

- The hybrid model integrated water balance equations with deep learning algorithm
- The proposed model presented the superiority and powerful simulation ability
- 5 The automatic parameterizing ability enhanced the model for cross-region simulation

1 1 Introduction

2 Groundwater plays an important role in geophysical and hydrological simulation (Cuthbert et al., 2019; Li et al., 2017; Zipper et al., 2018). Two major focuses related to groundwater 3 simulation are: 1) exploring the essence of hydrological cycle and its potential influences on 4 geophysical dynamic system, such as climate change (Ma et al., 2021a; Fan et al., 2013; Mohan et 5 al., 2018) and drought propagation (Ma et al., 2021b; Sadeghfam et al., 2018; Zhou et al., 2020); 6 7 and 2) exploring rational strategies including management, utilization and protection of groundwater resources under the impacts of human activities, such as pumping (Zipper et al., 2018) 8 9 and pollutant transport (Ossai et al., 2020; Tian et al., 2015). For groundwater-related simulations, 10 seeking accurate prediction of groundwater level (GWL) is an inevitable task. Due to the deepening understanding about the essence of geophysical dynamic system, many efforts have 11 been devoted into process-based hydrological models (Feng et al., 2020). In general, the successes 12 13 of traditional groundwater-related hydrological models relied heavily on detailed geological attributes and meteorological data. More specifically, referring to collected local hydrogeological 14 information, such as boundary conditions, land covers, and so on, a catchment is divided into 15 several hydrological units and water interaction between adjacent hydrological units is simulated 16 numerically (Mohan et al., 2018). Furthermore, strong mathematical derivation ability is required 17 for establishing groundwater-related models since the relevant governing equations are almost 18 two- or three-dimensional partial differential equations (PDEs). In general, process-based models 19 are the most suitable simulation strategies for catchments where the knowledge of flow path, 20 21 geological characteristics and boundary conditions are vital and extensively clear (Sahoo et al., 2017). Despite there are many successful cases of applying traditional process-based models to 22

simulate GWL in specific regions, it has been criticized as too complex, parametric, difficult and regional limited to use (Beven and Cloke, 2012; Clark et al., 2015a; Simone Fatichi et al., 2016; Tran et al., 2021). Moreover, it is still extremely difficult for process-based models to perfectly describe the water flow with nonlinear interactions, spatial heterogeneity, and time lags in real groundwater systems (Clark et al., 2015b). With the continuous development of computer science and deep learning algorithms (DLAs), data-driven hydrological models seem to be reliable alternatives for traditional process-based hydrological models.

In general, the main data sources for hydrological modeling originate from observations, 30 31 surveys, and lab experiments over years of systematic research (Chadalawada et al., 2020). Multiple variables, sources, and resolutions in space and time doom the complex and big-data 32 nature of hydrological data. With the aiming of exploring, extracting, and simulating the 33 information from the raw data, different types of DLAs emerged as times required. Among the DL 34 strategies, artificial neural networks (ANNs) and their optimized forms have become the 35 mainstream methods for establishing data-driven hydrological models. During the past decades, 36 applications of data-driven DL models have burgeoned in different fields of hydrology and their 37 performances have been widely recognized (Reichstein et al., 2019). For examples, commonly 38 39 recognized DLAs for hydrological timeseries modeling, such as streamflow and GWL forecasting, include: recurrent neural network (RNN) based models, for instance, long short-term memory 40 (LSTM) (Feng et al., 2020; Gao et al., 2020; Kratzert et al., 2019; Jiang et al., 2021) and gate 41 42 recurrent unit (GRU) (Cai et al., 2021; Zhang et al., 2021); convolutional neural networks (CNN) (Jiang et al., 2018; Jiang et al., 2019); genetic programming (GP) (Babovic, 2009; Babovic and 43 Keijzer, 2002; Chadalawada et al., 2017; Liu and Shi, 2019; Chadalawada et al., 2020); and etc. 44 Admittedly, benefited from the outstanding computational efficiency with diverse algorithm types, 45

DL models present excellent adaptability and versatility when serving as alternatives to traditional models. Comparing to the process-based models, DL models have advantages in terms of easy construction, computing speed, and data requirements (Reichstein et al., 2019). However, criticizes on DL models focus on the black-box essence, which means that DL models are more like computing tests without any physical meanings. As a result, DL models have not been psychologically accepted by hydrologists because such models cannot improve our understanding the essence of natural hydrogeological processes at present.

In recent years, application of artificial intelligence methods on hydrological modeling has 53 developed into a new and critical stage which integrating DLAs with dynamic geophysical 54 processes is expected to enhance the performance and generality of data-driven models 55 simultaneously (Frame et al., 2022; Jiang et al., 2020; Hoedt et al., 2021; Reichstein et al., 2019; 56 Zhao et al., 2019). Recent studies have demonstrated two strategies of integrating physical 57 constraints and DLAs: 1) using interpretation methods to demonstrate the prevalent DL 58 architectures for gaining scientific insights. For example, Jiang et al. (2022) explored the flood 59 inducing factors by analyzing inner works of the LSTM models. 2) embedding physical processes 60 into DL models to improve their awareness to systematic dynamic processes. Among these 61 62 research works, one strategy is adding physical informed equations into loss functions of deep learning models (Raissi et al., 2019). For example, Wang et al. (2020) proposed a Theory-guided 63 Neural Network (TgNN), which considered groundwater related factors (including governing 64 65 equations, boundary conditions, initial conditions, engineering controls and expert experience) as residual terms for loss function of the neural networks. This type of strategy is mostly suitable for 66 specific spatial problems as it could be considered as a powerful approximator for the PDEs. 67 Another strategy is adding physical constraints into the inner neural networks for forward 68

propagation. This type of strategy is majorly based on recurrent neural networks (RNN) for 69 timeseries simulations. For example, Kratzert et al. (2019) optimized the LSTM models with 70 regional entity awareness for streamflow simulation; Zhao et al. (2019) integrated Penman-71 Monteith equations into ANN models for simulation of evapotranspiration. Significantly, Niu et 72 al. (2019) demonstrated the relationship between the network architecture of RNN family and 73 numerical method, and theoretically supported the use of RNN to solve problems involving system 74 dynamics. Jiang et al. (2020) successfully embedded snowmelt process into DL models for 75 streamflow simulations. In general, most previous studies only focused on streamflow-related 76 77 tasks, applications of physically constrained DL models on the GWL simulation remains to be explored. 78

With the purpose of improving the performance of lumped DL model for GWL simulation 79 at catchment scale with limited geological observations, this study proposed a hybrid hydrological 80 model with water balance as physical constraints and DLA as cornerstone. The embedding of water 81 balance equations is theoretically supported by the algorithm of solving ordinary differential 82 equations (ODEs) with RNN (Jiang et al., 2020). The groundwater-related water balance equations 83 are summarized from GSFLOW model (Markstrom et al., 2008), EXP-HYDRO model (Patil & 84 85 Stieglitz, 2014), and TOPMODEL (Kirkby, 1975). Two prevalent DLAs for timeseries simulations, i.e., one dimensional convolutional neural network (1D-CNN) and gate recurrent unit 86 (GRU), are established for comparisons of simulating performance. Specifically, the hybrid model 87 88 consists of a self-designed RNN model (WB-Model) with wrapped water balance equations and a prevalent DL model (two-layer 1D-CNN model). Self-designed parameters of water balance 89 equations with physical meanings in WB-Model are determined during the iterative (training) 90 91 process of DLA, which strengthens the model's understanding of physical process and

hydrogeological characteristics of target catchments. For hybrid model, GWL simulated form WB-92 Model serves as an input for a 2-layer 1D-CNN model to give the final output. As a result, 93 enhanced simulation accuracy, robust, generalization ability and intelligence for inferring 94 characteristics of catchments are expected from the proposed hybrid model. Overall, this study 95 demonstrates that the novel hybrid model can garner the GWL-related physical knowledge in a 96 97 catchment vision if integrated with physical constraints properly, which makes the physical constrained DL model be more accurate, interpretable, feasible, advanced, and promising in terms 98 of GWL simulations for cross-region and less ungauged catchments. 99

100 **2 Data and methods**

101 2.1 Study area and multisource data

Continental United States (CONUS), as shown in subplot(a) of Fig. 1, is divided into 18 major 102 watersheds by U.S. Geological Survey (https://www.usgs.gov/), and each watershed contains 103 either drainage area of a major river or the combined drainage area of several rivers. Experimental 104 data collected in this study consist of two datasets from 10 major watersheds: New England Region 105 (01), Mid Atlantic Region (02), South Atlantic-Gulf Region (03), Great Lakes Region (04), Ohio 106 Region (05), Tennessee Region (06), Upper Mississippi Region (07), Lower Mississippi Region 107 108 (08), Missouri Region (10) and Arkansas-White-Red Region (11). The first dataset is from publicly available Catchment Attributes and Meteorology for Large-Sample Studies (CAMELS), which 109 contains two types of data to describe a specific catchment with minimal human disturbance: 1) 110 111 seven types of basin-averaged daily hydrometeorological timeseries data mostly recorded from 1980 to 2014 at hydrometeorological observation stations: precipitation (P), surface downward 112 solar radiation (SRAD), snow water equivalent (SWE), maximum temperature (T_{max}) , minimum 113 temperature (T_{min}) , near surface daily average vapor pressure (V_p) , and streamflow observations 114

at catchment outlet (SF); and 2) six types of averaged catchment attributes, i.e., topography and 115 location, climate indices, hydrological signatures, land cover characteristics, soil characteristics, 116 and geological characteristics (Addor et al., 2017; Newman et al., 2015). Each catchment in the 117 CAMELS dataset is represented by a hydrological unique code (HUC). The second dataset is the 118 freely available daily GWL data corresponding to catchments in the CAMELS dataset, which was 119 collected and compiled by U.S. Geological Survey. In this study, we selected the GWL data 120 following three principles simultaneously. First, the selected groundwater wells must provide daily 121 GWL data with 25 to 30 consecutive years between 1980 and 2014. Lower limit of 25 years was 122 123 set to ensure both enough monitoring wells that could meet the standard and a satisfying amount of data for the need of DL model. Meanwhile, upper limit of 30 years was set to control the 124 difference of data volume for DL model to ensure the simulation results were comparable. Second, 125 as shown in Fig. 1, GWL data must be collected from the monitoring wells closest to the 126 streamflow observation stations from the same catchment. This principle was set to make sure that 127 the hydrological factors in the CAMELS dataset were major driving forces of the GWL changes. 128 Despite we recognized that the boundaries of groundwater and streamflow catchments are not 129 always perfectly overlapped, among the dataset we collected, compared with the average 130 catchment area of 812 km^2 , we believed that the average distance of 20 km between groundwater 131 monitoring wells and streamflow stations was small enough to ensure their potential interactions. 132 133 Third, since the CAMELS dataset contained only catchments with minimal human impacts and 134 human impacts were not considered in this study, the groundwater data we collected should also 135 avoid urban and agricultural areas with frequent human activities. However, most of the wells we collected in the western CONUS, especially in California Region, presented obvious traces of 136 human activities such as pumping. The number of eligible wells in those regions is negligible. As 137

138 a result, 91 catchments from 10 major watersheds located in the central eastern CONUS were used



139 for this study (Fig. 1).

140



147 2.2.1 Gate recurrent unit

A general concept of RNNs refers to a class of ANNs with recurrent cells, which works on the principle of storing the output of a particular layer and feeding it back to the input in order to predict the optimate output of the layer. During the past decades, a variety of RNN-based DL models, including GRU and LSTM, were proposed following different strategies of designing the recurrent cell. To better introduce GRU model, we will firstly introduce the basic model of RNN family, i.e., simple RNN model. As shown in Fig. 2(a), a simple RNN (Rumelhart et al., 1986) recurses in the evolution direction of the sequence and all nodes are connected in a chain. The key point of a simple RNN model is the concept of hidden state (S_t), which stores information from previous time steps and delivers the previous features to predict the output (y_t). The specific algorithms are listed as follows:

$$S_t = \sigma_s (U_s x_t + W_s S_{t-1} + b_s) \tag{1}$$

$$y_t = \sigma_y \Big(W_y S_t + b_y \Big) \tag{2}$$

where σ is an activation function strategy, U and W are weight matrices, and b is a bias vector. 160 Theoretically, RNN can make use of all the information from former sequence, which makes it a 161 preferred strategy for timeseries analysis. However, due to the gradient vanishing and explosion 162 problems, a simple RNN would perform poorly when it comes to long sequence analysis because 163 outputs are likely to be only determined by several former steps. In order to solve these problems, 164 LSTM was firstly brought up for language processing with three gating controllers and two hidden 165 states (Hochreiter and Schmidhuber, 1997). To simplify the LSTM structure, GRU neural 166 networks (as shown in Fig. 2(b)) were proposed by Cho et al. (2014). GRU neural networks share 167 a similar chain structure with that of a simple RNN, but the internal operations inside the recurrent 168 cell are optimized for long short-term sequence simulation. The core algorithm of GRU is the two-169 170 gate controller: reset gate and update gate. The reset gate determines how to combine the new input information with the previous memory (which information will be stored or deleted), and the 171 update gate defines the amount of previous memory saved to the current state (which data will be 172 output from the current state). Computational workflow of GRU can be summarized as follows: 173

174
$$z_t = \sigma(W_z[x_t, h_{t-1}] + b_z)$$
(3)

175
$$r_t = \sigma(W_r[x_t, h_{t-1}] + b_r)$$
(4)

176
$$\widehat{h_t} = \tanh(W_h \cdot x_t + U_h \cdot (r_t \otimes h_{t-1}) + b_h) \tag{5}$$

$$h_t = (1 - z_t) \otimes h_{t-1} + z_t \otimes \widehat{h_t}$$
(6)

where W and b are weight matrices and bias vectors; z_t and r_t are update and rese gates' 178 activation vectors; \hat{h}_t is potential update vector; σ and tanh are sigmoid and hyperbolic tangent 179 functions. Recently, several studies have shown that GRU presents similar, if not better, simulation 180 performance with LSTM but a better computational efficiency due to a reduction of structural 181 complexity (Ayzel and Heistermann, 2021; Gao et al., 2020). Therefore, we applied GRU 182 networks as one scenario of GWL simulation. As shown in Fig. 2(d), referring to the modelling 183 184 strategy introduced by Cai et al. (2021), we adopted two GRU layers with a full connection layer and a dropout layer for the catchment-scale GWL simulation. 185

186

2.2.2 One dimensional convolutional neural network (1D-CNN)

187 Convolutional neural network was firstly proposed by LeCun et al. (1989), who integrated the back propagation algorithm and shared weights into convolutional neural layer. Normally, the 188 189 basic structure of CNN consists of input layer, convolutional layer, pooling layer, full connection 190 layer, and output layer. As shown in Fig. 2(c), 1D-CNN structure applied in this study adopts 191 convolution layer and pooling layer alternately. A convolution operation involves two steps: 1) 192 multiple learnable convolutional kernels (also known as filters) read the former layer, such as input 193 layer, by sliding on sequence matrix; and 2) the output features of the upper layer are convoluted with the convolution kernel, that is, the dot product operation is performed between the input term 194 and the convolution kernel, and then the results are sent to the activation function to obtain the 195 output features. The essence of convolutional layer is similar to RNN-based DLA, which is 196 extracting the most relevant features of input sequence for output sequence predicting. Comparing 197 with a fully connected layer, the 1D-CNN layers do not require the manipulation of lags 198 concatenation and decrease the memory resources allocated accordingly (Jiang et al., 2020). 199

200 Comparing with the RNN-based DL model, the algorithm of 1D-CNN is more concise, which may 201 make its simulation performance not be as good as that of RNN-based model, but its calculation 202 speed will be much faster (Jiang et al., 2020). To this point, 1D-CNN layers are chosen as one 203 section of the hybrid model for efficient testing.





Fig. 2. (a) Structure of a simple Recurrent Neural Network. (b) Structure of Gate

206 Recurrent Unit. (c) Structure of 1D-Convolutional Neural Network. (d) Workflow of

207 GRU-DL model in this study.

208

209

2.2.3 DLA for solving ODEs

In general, the simplest form of RNN-based model is that the input vector sequence $(x_1, x_2, ..., x_t)$ and output vector sequence $(y_1, y_2, ..., y_t)$ should satisfy the following recursive relationship:

$$y_t = f(y_{t-1}, x_t, t)$$
 (7)

As shown in Fig. 2, state (S_t) stores and propagates the information from last time step (S_{t-1}), which leads to the simplest RNN model is in fact a recursive computation model as follows:

216
$$y_t = \tanh(W_1 y_{t-1} + W_2 x_t + b)$$
 (8)

217 where W and b are weight matrices and bias vectors.

Research concerned with dynamics and dynamical systems is usually closely related to the solution of ODEs (Jiang et al., 2020). General ODEs have the following form:

220
$$\dot{y}(t) = f(y(t), t)$$
 (9)

221 Only a few ODEs can be solved analytically. In most cases, numerical method is regarded as the 222 first choice for solving ODEs. Euler method is a basic numerical method for solving initial value 223 problems of the first-order ODEs (Butcher, 2000). Implicit Euler method uses the first-order 224 backward difference quotient instead of differentiation, that is:

225
$$\dot{y} \approx \frac{y(k) - y(k-1)}{t(k) - t(k-1)} = \frac{y(k) - y(k-1)}{h}$$
 (10)

where h is time step. In this case, the differential equation becomes an implicit difference equation:

227
$$\begin{cases} y(k) = y(k-1) + hf(t(k), y(k)) \\ y(0) = y_0 \end{cases}$$
(11)

In the *k*-th iteration, y(k - 1) and t(k) are known values, y(k) is the unknown value to be obtained, and f(t(k), y(k)) is a function of the unknown value y(k) to be solved. Normally, y(k)can be solved by iterative solution of nonlinear equations (such as Newton's iterative method):

calculating the values of y(1), y(2), y(3), y(4), ..., y(n) in turns, and the discrete sequence is the 231 numerical solution of the differential equation. However, comparing Eq. 7 with Eq. 11, it can be 232 233 concluded that the Euler solution of ODE is theoretically a special case of RNN. The forward propagation of RNN corresponds to the Euler solution of ODE. Meanwhile, the backpropagation 234 process of RNN for determining the parameters brings out the advantage of applying RNN to solve 235 ODE (Tan et al., 2018). To be more specific, in actual problems, there is a type of problem called 236 "model inference", which is based on the known experimental data. The ultimate goal is to guess 237 the model (mechanism inference) that this batch of data conforms to. The approach to solve this 238 type of problem is roughly divided into two steps: the first step is to guess the form of the model, 239 and the second step is to determine the parameters of the model (Viana et al., 2021). In our case, 240 assuming that this batch of data can be described by an ODE whose form is already known, what 241 needs to be focused on is the estimation of the inner parameters. From the perspective of RNN, 242 the process of parameter determination corresponds to how the propagation mechanism of RNN 243 244 model is designed: forward propagation is to solve ODE (RNN prediction process), and backpropagation is to infer ODE parameters naturally (RNN training process). An interesting fact 245 is that, inferencing ODE parameters is a well-studied subject, but it is just one of the most basic 246 247 applications of RNN model.

As shown in Fig. 3, general idea of inserting physical constrained ODEs into the general concept of RNN structure is to replace the step function of simple RNN model with differentiable physical constrained equations. In our case, the proposed WB-Model is designed to encode geosystem dynamics as a neural network architecture with RNN backend. Within the self-designed RNN architecture, the connections (Eq. 1 and Eq. 2) between neurons in simple RNN structure, inputs, states, and outputs, are specified with the state-space representation in an explicit discrete

form of governing functions, which in our case are the water balance equations. Specifically, we 254 replaced the algorithms and the activation functions for calculating the states (S_t) in hidden layer 255 and output (y_t) by water balance ODEs with self-designed state variables, self-designed flux 256 variables and self-learnable parameters (see equations in Fig. 3(a) and Fig. 3(b), respectively). To 257 make it clearer, Fig. 3(c) presents the structure and unfolded network schematic graph of the WB-258 Model layer. Unlike simple RNN model, which possesses one state variable for each hidden layer, 259 WB-Model is equipped with five states variables (S_0, S_1, \dots, S_4) for each of the so-called hidden 260 261 layer. Each state variable represents a specific groundwater related hydrological reservoir. To calculate and update the state variables, nine flux variables are introduced. Each flux variable 262 represents a flow process. Flux variables are like the 'bridges' connecting the state variables in 263 order. To make the flux variable accurate and meaningful, 16 self-learnable parameters (θ_P), 264 consisting of θ_P^f and θ_P^g , with physical meanings are introduced. Each parameter is determined by 265 the gradient descent algorithm wrapped in the ordinary RNN model. Relevant descriptions about 266 the water balance equations, self-designed state variables, self-designed flux variables, and self-267 learnable parameters applied in this WB-Model are introduced in detail in subsection 2.3. In 268 summary, we redefined the step function of ordinary RNN to solve the target ODEs by introducing 269 self-designed state variables, self-designed flux variables, and self-learnable parameters with 270 271 physical meanings.





273

277

Fig. 3. Architectures of the simple RNN layer and the proposed WB-Model layer.

274 2.3 Groundwater-related water balance equations

The essence principle for process-based models is water balance, a form of mass conservation law in terms of hydrological research (Herrmann et al., 2016), with a generic equation as:

$$\frac{dS}{dt} = I(t) - O(t) \tag{12}$$

where S is the water storage unit, I(t) is the inflow to the unit at time t and O(t) is the outflow from 278 279 the unit at time t. Traditional process-based hydrological model usually divides the catchments into segments (hydrological units) and calculates the flow between adjacent segments based on 280 water balance equations and hydrogeological characters of segments. In addition to spatial and 281 temporal resolutions, existing process-based hydrological models majorly differ in mathematical 282 algorithm of I(t) and O(t) (Jiang et al., 2020). In this subsection, we summarized water balance 283 equations related to GWL from different process-based models including GSFLOW model 284 (Markstrom et al., 2008), EXP-HYDRO model (Patil & Stieglitz, 2014) and TOPMODEL (Kirkby, 285

1975). These processes were integrated into DL by replacing the states and learnable parameters
with water balance equations for RNN model. We named the model and equations wrapped in the
RNN model as WB-Model, and related processes and equations are listed below.

289

306

2.3.1 Canopy interception

As shown in Fig. 4, interception of precipitation by plant canopy is computed during a time 290 step as a function of plant-cover density and the storage available on the predominant plant cover 291 type. The variability of the CAMELS dataset makes it possible to demonstrate the plant canopy 292 interception process since the landcover related data is included (Dunkerley and Booth, 1999). 293 Precipitation that reaches the ground is calculated by the sum of throughfall and precipitation on 294 the catchment not covered by plants. Canopy interception is closely related to the vegetation type 295 and local climate condition (Tao et al., 2020; Trinh and Chui, 2013). In this subsection, we set a 296 reservoir (S_0) for the storage change of canopy interception. Two critical learnable parameters, the 297 maximum capacity of canopy interception (SC_{max}) and canopy interception coefficient related to 298 299 plant type (K_c) , were introduced for physical RNN model. Since CONUS locates in the Northern Hemisphere, we determined the seasons according to the averaged day-length: if the day-length is 300 longer than 0.5, we set a range with higher mean value of K_c and SC_{max} because vegetation 301 flourishes in autumn and summer; conversely, if day-length is shorter than 0.5, we set a range with 302 smaller mean value of K_c and SC_{max} for winter and spring. The equations (Leavesley et al., 1983) 303 304 can be summarized as:

305
$$P_{int} = \begin{cases} K_c \times D_c \times A, 0 \le S_0 \le SC_{max} \\ 0, S_0 < 0 \\ P_{max}, S_0 > SC_{max} \end{cases}$$
(13)

$$\frac{dS_0}{d_t} = P_{int} \tag{14}$$

307
$$\begin{cases} K_c = K_{cs}; \ SC_{max} = SC_{smax}, when \ DayL \ge 0.5\\ K_c = K_{cw}; \ SC_{max} = SC_{Wmax}, when \ DayL < 0.5 \end{cases}$$
(15)

308 where D_c is vegetation coverage, which reflects the spatial coverage of vegetation. A is catchment

area. Values of D_c and A are from the CAMELS dataset.



310

Fig. 4. The main hydrological processes considered in the proposed model.

312

2.3.2 Snow melt

Snow melt process, a critical factor for hydrological simulation especially in cold regions, has been proven to be a key point for streamflow simulations with DL methods (Broxton et al., 2019). Inspired by the methods of wrapping EXP-HYDRO model (Patil and Stieglitz, 2014) into P-RNN model (Jiang et al., 2020), the snow melt process is implanted into the hybrid model as follows.

First, the precipitation reaches the ground will be divided into snowfall (P_s) and rainfall (P_r), which is controlled by daily temperature (T) and threshold temperature of snowfall (T_{min}):

320
$$\begin{cases} P_r = P - P_{int}; P_s = 0, when T > T_{min} \\ P_s = P - P_{int}; P_r = 0, when T \le T_{min} \end{cases}$$
(16)

321 Then, water from snow melt was calculated as:

322
$$M = \begin{cases} \min\{S_1, D_f \times (T - T_{max})\}, & if S_1 > 0 \text{ and } T > T_{max} \\ 0, & otherwise \end{cases}$$
(17)

where D_f is thermal degree-day factor, T_{max} is temperature threshold of snowmelt, and S_1 was the snow reservoir, which can be expressed as:

$$\frac{dS_1}{d_t} = P_s - M \tag{18}$$

326

325

2.3.3 Soil retention and groundwater flow

Describing the dynamic system of water movement in soil has always been a complex but critical object for groundwater simulation (Banwart et al., 2019; Cai et al., 2021). Indeed, soil structure affects the all-round local hydrological response (Fatichi et al., 2020). What we focused on was the water distribution in soil that meets the catchment-scale water balance requirements. Again, this part is referred to the soil structures introduced in GSFLOW model. Specifically, we only considered the downward vertical flow process between preferential reservoir, capillary reservoir, and gravity reservoir, with the aim of meeting the catchment-scale water balance.

334

2.3.3.1 Preferential reservoir

When precipitation and snow melt reach the ground, a fraction of infiltration is apportioned to the preferential-flow reservoir to account for fast interflow through large openings in the soil zone near land surface, while the rest of the water, which is generated when the precipitation rate exceeds the infiltration rate of the soil that may not be saturated will be classified as Hortonian flow (Horton, 1933). Eventually, both Hortonian flow and preferential flow contribute to the streamflow of catchment outlet. Referring to the TOPMODEL and P-RNN model, water flow in the preferential-flow bucket can be calculated as:

342
$$Q_{hor} = (P_r + M) \times (1 - R_{in})$$
 (19)

343
$$Q_{pref} = \begin{cases} 0, & S_2 < 0\\ (P_r + M) \times R_{in} \times R_{pr} \times e^{-f \cdot (S_{pmax} - S_2)}, 0 \le S_2 \le S_{pmax} \\ Q_{pmax}, & S_2 > S_{pmax} \end{cases}$$
(20)

where R_{in} is the infiltration rate; R_{pr} is the coefficient of fast interflow; f is decay factor; S_2 is preferential reservoir; S_{pmax} is the storage capacity of preferential bucket; and Q_{pmax} is the maximum preferential flow. The values of the parameters are learned during iterative process of self-designed RNN model, which can be expressed as:

$$\frac{dS_2}{d_t} = (P_r + M) \times R_{in} \times R_{pr}$$
(21)

349

355

356

2.3.3.2 Capillary reservoir

The capillary reservoir represents water held in the soil by capillary forces between the wilting and field-capacity thresholds. Water is removed from the reservoir by evapotranspiration (Markstrom et al., 2008). Referred to the EXP-HYDRO model and P-RNN model, the PET (Potential Evapotranspiration) is estimated by Hamon's formulation. Therefore, the calculation can be concluded as follows:

$$Q_{cap} = (P_r + M) \times R_{in} \times (1 - R_{pr})$$
(22)

$$\frac{dS_3}{dt} = Q_{cap} - ET \tag{23}$$

357
$$ET = \begin{cases} 0, S_3 < S_{cmin} \\ PET \times \left(\frac{S_3}{S_{cmax}}\right), 0 < S_3 < S_{cmax} \\ PET, S_3 > S_{cmax} \end{cases}$$
(24)

358
$$PET = 29.8 \times L_{day} \times \frac{0.611 \cdot e^{17.3 \cdot \frac{T}{T + 237.3}}}{T + 237.3}$$
(25)

where Q_{cap} is the waterflow into capillary reservoir, S_3 is capillary reservoir, S_{cmin} and S_{cmax} are the minimum and maximum storage capacities of capillary reservoir, L_{day} is day length and *T* is 361 local temperature. The main function of capillary reservoir is to calculate the amount of 362 evapotranspiration. The maximum storage capacity represents the field capacity, and the minimum 363 storage capacity represents the minimum storage capacity held by vegetation.

364

375

2.3.3.3 Gravity reservoir

The gravity reservoir represents water in the soil zone between field-capacity and saturation 365 thresholds. This reservoir was developed to provide gravity drainage from the soil zone to the 366 unsaturated zone, which will eventually discharge to the groundwater. According to GSFLOW 367 model, slow interflow from the gravity reservoirs represents the perching of water in the soil zone 368 above the water table that can occur because of mineralization near the bottom of the soil zone or 369 when soil develops over fine-grained material. Slow interflow can occur when the water content 370 in the soil zone exceeds the field-capacity threshold. Slow interflow is developed from continuity 371 and an empirical equation (Leavesley et al., 1983). For the lumped hydrological model, we added 372 a decay process for slow interflow referring to TOPMODEL, which can be written as: 373

374
$$Q_{slow} = \begin{cases} 0, & S_4 < 0\\ Q_{gra} \times k_l + Q_{gra}^2 \times k_n \times e^{-f \cdot (S_{gmax} - S_4)}, & 0 \le S_4 \le S_{gmax}\\ Q_{gmax}, & S_4 > S_{gmax} \end{cases}$$
(26)

$$Q_{gra} = Q_{cap} - ET \tag{27}$$

$$\frac{dS_4}{d_t} = Q_{slow} \tag{28}$$

where Q_{gra} is the water flows into gravity reservoir; S_4 is gravity reservoir; S_{gmax} is maximum storage capacity of gravity reservoir; k_l is linear coefficient of slow interflow; k_n is non-linear coefficient of slow interflow; and Q_{gmax} is the maximum slow interflow.

Finally, the outputs of the water balance model are expressed as:

$$Q_{stream} = Q_{hor} + Q_{pref} + Q_{slow}$$
(29)

$$Q_{gw} = Q_{gra} - Q_{slow} \tag{30}$$



386

To better illustrate the way of wrapping above equations into WB-Model, we use snow melt process in subsection 2.3.2 as an example. The precipitation that reaches the ground (flux from canopy interception reservoir) and temperature in input (x_t) will be used for initial input for this section, the snow reservoir. Firstly, we set a temperature threshold T_{min} ranging from -3°C to 0°C to determine whether precipitation reaching the ground is in form of rainfall or snow (Eq. 16).

Secondly, we set a temperature threshold T_{max} ranging from 0°C to 3°C to determine when the 392 snow starts to melt. A thermal degree-day factor D_f ranging from 0 to 5 (mm/day/°C) was set to 393 calculate the snow melt (Eq. 17). The state variable (S_1^t) , which represent the snow storage, is 394 updated by Eq. 18 as $S_1^t = S_1^{t-1} + P_s^t - M^t$. In this case, T_{min} , T_{max} and D_f will be determined 395 through the training process. Above process is inserted into the neural network by defining S_1 , 396 rainfall, snow and melt with Eq. 16 to Eq. 18. Further, the flux out of snow reservoir (S_1) , which 397 is the rainfall and snow melt, will serve as the input for preferential reservoir in subsection 2.3.3.1 398 to further update S_2 and calculate other fluxes. With similar process, state variables will be updated 399 in order from S_0 to S_4 indicating water flows from precipitation to groundwater. As a result, after 400 updating the five states and nine flux variables, the preliminary groundwater fluctuation would be 401 402 calculated as Eq. 30. Pseudocode of the ODE-RNN layer for GWL simulation is presented in Table 1. Overall, the process-based groundwater model involves three input daily variables from the 403 CAMELS dataset (Precipitation, averaged Temperature and Day-length), five states variables (S_0 , 404 S_1 , S_2 , S_3 and S_4), nine flux variables (P_{fall} , P_s , P_r , M, Q_{cap} , Q_{pref} , Q_{hor} , Q_{gra} and Q_{slow}), and 16 405 learnable parameters, which controls the hydrological behaviors (SC_{max} , K_c , T_{min} , T_{max} , D_f , R_{in} , 406 $R_{pr}, S_{pmax}, Q_{pmax}, f, S_{cmax}, S_{min}, k_l, k_n, S_{gmax}, Q_{gmax}$). Again, the proposed WB-Model is a 407 spatially lumped DL model with physical constrains, which adheres strictly to the law of water 408 balance. 409

410 2.4 Integrated hybrid framework for GWL simulation

Inspired by the outstanding work of P-RNN integrated hybrid streamflow simulation model (Jiang et al., 2020) and previous work of exploring the usage of DL model for large-scale GWL simulation (Cai et al., 2021), this study proposed a hydrology-aware DL architecture for GWL simulation. As presented in Fig. 5, there were two pipelines for the hybrid model. The first pipeline

was the water balance based RNN model (WB-Model), which wrapped catchments attributes and 415 water balance equations into DLA. The forcing variables for this pipeline were precipitation, day-416 length, and temperature. The parameters within WB-Model were trained separately to present a 417 preliminary GWL, which was furtherly served as one of the input variables for the second pipeline. 418 The second pipeline was a sequence-to-sequence DL model, which maps the meteorological 419 420 sequences to GWL sequence. The input variables of the second pipeline included the preliminary GWL simulated by WB-Model and the other six dynamic forcings provided by the CAMELS 421 dataset. The objective of WB-Model was to identify geo-hydrological information from 422 423 observations of the external world and reorganize them in form of providing simulated preliminary GWL. In other words, WB-Model was used to facilitate the prevalent DL model. The core function 424 of WB-Model was to combine water balance equations and DL algorithm to capture the fluctuation 425 characteristics of GWL that response to driven factors. As a result, the preliminary GWL simulated 426 by WB-Model was served as an important input forcing to strengthen the training process of hybrid 427 model. Considering the ability of handling lagged effect from hydrological signals (Feng et al., 428 2020) and a faster computing speed than RNN-based model (Jiang et al., 2020), a two-layer 1D-429 CNN model was chosen as the main DLA for simulating the GWL. After a large number of 430 431 preliminary experiments considering the trade-off between increasing modelling accuracy and reducing the complexity of the model structure, the first layer applied 8 kernel filters, each with a 432 length of 15. The length of 15 symbolized that the influence on the current hydrological response 433 434 could be traced back to 15 days ago. From a hydrological viewpoint, tracing the influences on the current hydrological response back to 10 days ago is a common strategy in data-driven 435 hydrological models and has been proven to be successful for streamflow simulation with similar 436 437 method (Jiang et al., 2020). Since our simulation target is GWL for phreatic aquifer, the

hydrological response time will be a little longer than that of streamflow simulation. The set of 15-438 day hydrological response time was based on a large number of preliminary experiments. 439 Nevertheless, we cannot ensure that 15-day hydrological response time is the most accurate setting 440 for each specific catchment. However, this setting is overall feasible for GWL simulations of these 441 91 catchments and is in line with common strategy in data-driven hydrological models. The second 442 layer used 1 convolution kernel filter for analyzing the output of the first layer and providing the 443 final results. Zero-padding strategy was adopted for ensuring the output had the same sequence 444 length as the input sequence after filtering by kernel filters. 445



446

447 Fig. 5. The proposed generic architecture that explicitly embedding water balance
448 constrained dynamic behaviors into DL models.

The performance of the model was evaluated with Nash-Sutcliffe coefficient of Efficiency
(NSE) value (Nash and Sutcliffe, 1970) as follows:

451
$$NSE = 1 - \frac{\sum_{i=1}^{n} (Q_s - Q_o)^2}{\sum_{i=1}^{n} (Q_o - \overline{Q_o})^2}$$
(31)

where Q_s , Q_o and $\overline{Q_o}$ are the simulated, observed and mean observed GWL. NSE is widely used 452 for model assessment, which ranges from $-\infty$ to 1. The closer NSE is to 1, the better the simulation 453 result is. The NSE value was also set as loss function for all the DL models in this study. 80% of 454 each timeseries was divided as training set while the rest 20% of the data was used as testing set. 455 The epochs number was set as 400 and the learning rate was 0.01. The number of parameters in 456 hybrid model is 593 including 16 self-designed parameters in WB-Model, 568 parameters in first 457 1-D CNN layer and 9 parameters in the second 1-D CNN layer. All the models and equations 458 introduced in this study were coded with Python 3.6 with Keras as coding API and Tensorflow 1.4 459 as the DL backend with no GPU requirement. 460

Generation ability (GA) is an important evaluation criterion for data-driven model (Chen et al., 2020). It reflects the ability of DL model, which is well trained, to digest new data and make accurate predictions. The calculation method is as follows:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (O_s - O_o)^2}$$
(32)

$$GA = \frac{RMSE_{pred}}{RMSE_{train}}$$
(33)

464

where O_s and O_o are the simulated and observed GWL; *RMSE* refers to Root Mean Square Error, which measures the average magnitude of the error between model simulations and observations; and *GA* is calculated as the ratio of the *RMSE* values between predicting process and training

469	process (Yoon et al., 2011). If the data-driven model simulates GWL perfectly, the GA values will
470	be unity. Otherwise, if the model is over trained or overfitting, the GA values will exceed unity. In
471	this case, the model would not give accurate predictions even though it is able to provide accurate
472	fittings for the training data. The GA values will be less than unity if the model is under trained,
473	which is also called as underfitting.

- 474 **3 Results and discussion**
- 475 3.1 Comparisons of GWL simulation results

476 3.1.1 Comparisons of model performances in specific catchments

As representatives of the 91 catchments, three catchments (HUC: 04216418, HUC: 03015500 477 and HUC: 03026500) were chosen for presenting, describing, and analyzing the simulations in 478 detail. Fig. 6 presented the GWL forecasting performances with four independent DL strategies, 479 from top to bottom of the subplots are Hybrid model, 1D-CNN model, WB-Model (the self-480 designed RNN model), and GRU model. WB-Model integrated water balance constrains into 481 RNN, which served as a major optimization for hybrid DL model while the 1D-CNN model and 482 GRU model are pure DL model. Consequently, what concerned us the most was how much the 483 simulation performance of hybrid model was improved compared with the 1D-CNN model and 484 485 GRU model. The reason why we chose these three catchments as representatives was that the simulation performances of the 91 catchments could be roughly divided into three categories, 486 which corresponded to the three catchments shown in Fig. 6. 487





490

489 **Fig. 6**

Fig. 6. Representatives of simulation results by four DL strategies introduced in this study.

The simulation results of HUC: 04216418 represented the catchments that with an overall 491 492 simulation improvement when combining water balance constrains into common DL models. As shown in the first column of Fig. 6, the pure DL strategies, 1D-CNN (NSE=0.238) and GRU 493 (NSE=0.278) model, had similar simulations effect while the hybrid model (NSE=0.553) 494 outperformed significantly than the pure DL models. Such improvement was majorly benefit from 495 the accuracy of simulating the GWL fluctuance pattern by the WB-Model (NSE=-0.377). 496 Admittedly, based on the evaluation criteria of hydrological model, the performance of WB-Model 497 itself were not satisfied enough for an alternative of traditional process-based model. However, 498 although there were certain differences in specific values, the variation trend and fluctuation 499 patterns of the simulation results were consistent with the observations. As a result, the preliminary 500 result from WB-Model, served as an important input forcing, strengthened the training process of 501

hybrid model, which lead to a higher accuracy than pure DL models. More specifically, in plain language, common DLAs possessed the ability to filter out critical information that mostly related to the simulation target by adjusting the intrinsic parameters, such as weights and bias. In this case, comparing to the other two representatives, preliminary GWL simulated by WB-Model provided the most useful sequence for mapping the observations because the fluctuance pattern simulated by WB-Model was most similar to that of the observed GWL.

The second column of Fig. 6, HUC: 03015500, presented the representative of catchments 508 with another type of simulation performance. In general, despite the NSE improvements may not 509 510 be as obvious as the catchments in the first category, hybrid models (NSE=0.606) outperformed the pure DL models (NSE_CNN=0.391 and NSE_GRU=0.374) in these catchments for better 511 simulating the GWL fluctuance under extreme conditions. Comparing to the first category, the 512 overall simulation performances of WB-Model in this category were slightly inferior, which led 513 to a less obvious promotions of NSE values. As the simulation of catchment HUC: 03015500 514 illustrated, there was an obvious increasement of GWL during 2011 to 2012, which was resulted 515 from a large amount of snow melting and rainfall as precipitation. As shown in the 1D-CNN model 516 and GRU model, the pure DL models failed to capture such precipitation pattern because 517 518 precipitation characteristics, especially the extreme precipitation conditions, had certain probability of being diluted in a variety of input variables. During the training process, DL model 519 would consider the extreme precipitation events as outliers especially when the other inputs were 520 521 sufficient to establish a good liner or nonlinear mapping relationship with the object. This could be inferred from the fact that pure DL models had already presented good overall simulation 522 performances. Under these circumstances, WB-Model captured the pattern of extreme 523 524 precipitation events when providing the preliminary GWL, which emphasized the importance of the extreme points and reinforced the hybrid model to learn such patterns. Consequently, the simulation results of hybrid model were intuitively improved because of better matchings of extreme GWL.

52	8

529

Table 2. Summary of the simulation results of testing set by hybrid, 1D-CNN, and GRUmodel across 91 catchments in the 10 major watersheds.

Watersheds	Basin	1D-CNN		GRU		Hybrid Model		Improvement	
(Region)	count	Median	Mean	Median	Mean	Median	Mean	Mean	Count
01 New England	5	0.275	0.282	0.277	0.291	0.278	0.330	-0.008	3
02 Mid Atlantic	20	0.299	0.305	0.313	0.326	0.393	0.405	0.125	17
03 South Atlantic-Gulf	15	0.256	0.326	0.260	0.313	0.356	0.361	0.254	13
04 Great Lakes	5	0.365	0.348	0.373	0.347	0.528	0.578	0.231	5
05 Ohio	16	0.329	0.349	0.338	0.358	0.388	0.459	0.036	10
06 Tennessee	6	0.243	0.210	0.300	0.293	0.334	0.334	0.106	6
07 Upper Mississippi	8	0.298	0.354	0.443	0.275	0.475	0.475	0.097	7
08 Lower Mississippi	3	0.603	0.533	0.663	0.559	0.531	0.531	-0.013	2
10 Missouri	7	0.192	0.288	0.434	0.344	0.224	0.224	0.127	4
11 Arkansas-White-Red	6	0.314	0.401	0.432	0.548	0.409	0.409	0.244	4
ALL	91	0.279	0.253	0.313	0.278	0.408	0.401	0.129	71

530

The last column of Fig. 6, HUC: 03026500, was presented as a representative of catchments 531 532 with negative effects when integrating WB-Model. Intuitively, preliminary GWL simulated by 533 WB-Model failed to match the fluctuance of observed GWL (NSE=-6.759). Conversely, the pure DL models showed satisfying simulation results (NSE_CNN=0.626 and NSE_GRU=0.686). The 534 535 reduction of hybrid model performance (NSE=0.542) resulted from the failure of WB-Model. Specifically, for catchment HUC: 03026500, the GWL simulated by hybrid model showed a larger 536 537 fluctuation range than observed GWL, which was consistent with simulation from WB-Model. 538 There were three potential reasons: 1) the groundwater observation well and the corresponding

streamflow observation station located in different drainage areas. As mentioned in subsection 2.1, 539 although the groundwater data we collected were from the observation well that was closest to the 540 streamflow observation station at the catchment outlet, we cannot guarantee that the GWL was 541 driven by the precipitation from the streamflow catchment because the catchments of groundwater 542 and surface flow might not be geographically coincident. In this case, the WB-Model driven by 543 544 the precipitation from the CAMELS dataset would inevitably lead to failed GWL simulation results because the fluctuation of GWL was driven by a different precipitation input. 2) The water 545 balance algorithm was not applicable for local GWL simulation. The essences of the WB-Model 546 and hybrid model were physical constrained DL models for lumped GWL simulation, which meant 547 the detailed geological information related to groundwater flux were not collected or considered 548 in these models. This would lead to simulation mismatches if discharge and recharge of local 549 groundwater were not majorly from the precipitation but from other sources such as confined 550 aquifer or lakes nearby. In this case, the WB-Model would not provide a rather precise simulation 551 result even if the precipitation was one of the driven forces for the local groundwater fluctuance. 552 3) The ranges of learnable parameters were not set properly. This was also resulted from the lack 553 of geophysical information to determine the range of self-designed parameters especially the soil 554 555 related parameters. Overall, the influence from the third problem could be minimized by multiple tests while the first and second problems might be the most likely causes of the failures of WB-556 Model and hybrid model for GWL simulations in several catchments. 557

558

559

560

561

3.1.2 Comparisons of overall accuracy between hybrid model and pure DL models As shown in Table 2, 1D-CNN, GRU and water balance constrained hybrid model were applied to simulate GWL of 91 catchments located in 10 major watersheds of the middle eastern

CONUS. More specifically, the median NSE values of the simulation results by using GRU model

was 0.313 while the average NSE was 0.278, which was consistent with the NSE range of former 562 simulations in similar areas by applying GRU algorithm (Cai et al., 2021). As an advanced DLA 563 focused on timeseries simulation, GRU model served as a group of controlled experiments in this 564 study for physical constrained hybrid model. In the meanwhile, the median NSE value of 565 simulation by using 1D-CNN model was 0.297 and the average NSE was 0.253. In views of the 566 NSE values, the simulation performance of the two pure DL models were not much different, but 567 due to the complexity of model structure and the difference in emphasis of simulation (GRU model 568 focuses on timeseries simulation and CNN model focuses on image processing), the overall result 569 570 of 1D-CNN model simulation would be slightly inferior to that of GRU model for GWL simulations. Benefited from the WB-Model, the overall simulation performance was significantly 571 better than that of 1D-CNN and GRU model: the median NSE was 0.408 and the average value 572 was 0.401. The NSE distribution was shown in Fig. 7(b). 573

Since the final algorithm of hybrid model was a two-layers 1D-CNN model, comparing to 574 the pure 1D-CNN model, the preliminary GWL simulated by the WB-Model was the key factor 575 for improving the performances of hybrid model. With the aim of testing the adaptability of hybrid 576 model, we counted the number of catchments where the hybrid model outperformed the 1D-CNN 577 578 model for GWL simulation. As shown in the last two columns of Table 2, comparing with the 1D-CNN model, 71 out of 91 catchments presented simulation improvement when applying the hybrid 579 model. Furthermore, for the 10 major watersheds, comparing to the 1D-CNN model, average NSE 580 581 improvements while using hybrid model were positive except for the New England and Lower Mississippi regions. However, the number of catchments used for experiments was so small that 582 the average NSE improvement would be seriously affected by a single catchment with a major 583 584 failure of hybrid model. Moreover, more than half of the catchments in New England and Lower Mississippi region presented positive improvement when using hybrid model than using 1D-CNN model. Optimizing effect of hybrid model was more intuitively explained in Figs. 7(a) and 7(c). Distribution of catchments where hybrid model outperformed the 1D-CNN model were shown in Fig. 7(d).





Fig. 7. (a) Histogram for comparison of simulation result from hybrid, 1D-CNN and
GRU model; (b) Distribution of NSE values for simulation results of hybrid model; (c)
Box plot for comparison of simulation result from hybrid, 1D-CNN and GRU model; (d)

593 Distribution of catchments where hybrid model outperformed 1D-CNN model.

Lumped hybrid DL model for streamflow simulation based on snow melting processes had been tested to be successful for catchments in high latitudes (Jiang et al., 2020). In this study, in addition to the snow melting process, groundwater-related soil water flow processes were also properly embedded into the WB-Model. As a result, the hybrid model outperformed the pure DL model not only in the snow-domain areas but also basins in middle and low latitudes. Preliminary GWL simulated by WB-Model improved the noise tolerance of hybrid model to groundwaterirrelevant input features and strengthens the learning ability of the hybrid model to GWL changes. Therefore, comparing to the pure DL models, hybrid model possessed better adaptability and understanding of different hydrometeorological and geophysical conditions with the constrains of water balance equations. In conclusion, comparing to prevalent pure DL models, hybrid model exhibited significantly higher overall accuracy for GWL simulations.

605 3.2 Comparisons of model generalization ability

The generalization ability was evaluated by Eq. 33. If the model focused on the training 606 process rather than a general system, the GA value would be higher. This meant the model will 607 have overfitting problems. The higher the GA values were, the weaker the generalization ability 608 the model processed. Fig. 8 presented the distributions of the GA values for the four DL models 609 in form of boxplot. The average GA values of hybrid, WB-Model, CNN and GRU model were 610 1.269, 1.124, 1.431 and 1.396, respectively. This indicated the developed hybrid model presented 611 a better generalization ability than that of CNN and GRU model and WB-Model presented the best 612 overall generalization ability among the four DL models. It should be noted that WB-Model was 613 614 a group of ODE equations coded with the language used in recurrent neural networks. This meant that WB-Model was essentially a physic guiding model. For WB-Model, the water balance 615 equations, served as a backbone of training and predicting stages, were universal theorems. This 616 617 was the main reason why WB-Model presented the best generalization ability. Furthermore, the result was consistent with conclusions of the previous study, which proved that physic based 618 numerical models tended to present better generalization abilities than that of pure machine 619

learning models (Chen et al., 2020). In conclusion, benefited from WB-Model, hybrid model
presented superiority than pure DL models in terms of generalization ability.



622

Fig. 8. Comparison of the GA values between Hybrid, WB-Model, CNN and GRUmodel.

625

3.3 Robustness of the proposed model

In general, in addition to the advanced core algorithms, an important factor affecting the effect 626 of DL model was the amount of sample data as well as the potential relationship between input 627 628 data and target data. An excellent DL models should be equipped with the ability to ensure its own 629 stability while satisfying the accurate simulation performance while reducing the effect from the 630 interference information as much as possible (Su et al., 2018). To demonstrate the reliability and stability of the proposed model, we chose the catchment HUC: 03182500 for scenario tests. We 631 632 evaluated the GWL simulation results of hybrid model with six combinations of input features. 633 More input features signified a larger amount of input data, which could be either useful information or noisy for simulations of targets. In this study, since the WB-Model needed three 634 inputs (precipitation, temperature, and day length) for the self-designed RNN model, these three 635

variables were the minimum input requirement. As shown in Fig. 9, with the increasing number 636 of input features, NSE values of the hybrid model also presented an increasing trend, from 0.564 637 to 0.650. It was worth mentioning that although the NSE values of the three inputs strategy was 638 the smallest, the value of NSE=0.564 was a fine result for GWL simulations. Meanwhile, in 639 perspective of the NSE values, the difference of simulation performances between the six scenarios 640 was relatively small, which proved that the hybrid model processing a good stability. Furthermore, 641 it could be inferred that, comparing to SRAD and vapor, applying streamflow data as one of the 642 input features could better improve the accuracy of the hybrid model, which implied that the 643 feature of an important hydrological phenomenon, flow exchange between streamflow and 644 groundwater through baseflow, was captured from the perspective of hybrid DL model. This 645 reflected that DL models, which had always been criticized as black-box models without any 646 physical meanings, could also express some natural processed implicitly (Jiang et al., 2022). 647



Fig. 9. Comparisons of hybrid model performances with different input strategies.

650

648

649

Fig. 10 presented the comparisons of the stability of physical constrained hybrid model and 1D-CNN model as representative of pure DL model through the NSE values. We adopted the data from catchment HUC: 03182500 for testing experiments to compare the performance of the model under same input feature strategies. The result showed that hybrid model outperformed the pure DL model for every input strategy and the NSE value of hybrid model presented a smaller fluctuation range (from 0.564 to 0.651) than that of pure DL model (from 0.417 to 0.517), which implied the hybrid model performed more stably than pure DL model.

The reason for this phenomenon was that the WB-Model in the hybrid model could adjust the 658 values of the internal learnable parameters according to the regional characteristics, so that the 659 preliminarily simulated GWL could reach a certain accuracy. Then, as an input, the preliminarily 660 simulated GWL could help the DL model better understand the change characteristics of the GWL 661 to be simulated, so that the performance the hybrid model was obviously better than that of the 662 pure DL model. In conclusion, the physical constrained mechanism made the stability of the hybrid 663 model better than the pure DL model, which was another outstanding advantage in addition to the 664 higher accuracy of the simulation results of the hybrid model. 665



Fig. 10. Comparison of simulation performances between hybrid and 1D-CNN model withdifferent input strategies.

669

670

666

3.4 A rethink of relationships between simple RNN, WB-Model and GRU.

It is worth clarifying that simple RNN, WB-Model and GRU are based on a general concept 671 of RNN, which is a family of neural networks with recurrent cells. In general, simple RNN model 672 does not possesses the long-term memory ability because of the potential gradient vanishing and 673 674 explosion problem when updating the hidden states during the training process. The main difference between simple RNNs and GRUs (or LSTMs) is the latter introduces a so-called gating 675 mechanism into the recurrent cell to control the information flow, to address the gradient vanishing 676 and explosion problem. Likewise, in our WB-Model, we also modified the recurrent cells of the 677 general RNN models, where the nodes were connected using physical equations rather than the 678 usual perceptions. More specifically, GRU introduces reset gate and update gate to determine 679

whether the previous information involved in the cell states is added or discarded. In comparison, 680 in the WB-Model, the information flow is controlled using the thresholds as in hydrological 681 models. For example, if the temperature is larger than a threshold, the rainfall information will 682 become a part of the soil state; otherwise, it will be added to the snow state. Further, although the 683 specific values of these parameters are determined during the training process of DL model, the 684 essence of the WB-Model is a group of ODE equations (presented by TOMODEL, GSFLOW, and 685 other hydrological models) coded with the language used in neural networks, which makes the 686 WB-Model naturally have long-term memory in the form of hydrological states (e.g., snowpack 687 and preferential flow pack). In conclusion, the GRU and WB-Model model have taken distinct 688 ways to reform the recurrent cells of general RNN models, which could be considered as two 689 optimizing strategies of simple RNN model. 690

691

3.5 Limitations and potential improvements

692 Comparing with pure DL models, the proposed hybrid model for cross-region GWL 693 simulations at catchment scale is significantly improved in terms of the accuracy, generality, 694 robustness and physical meaning of models, but there are some obvious limitations and potential 695 improvements for the hybrid model.

The first limitation is that the process of adjusting self-designed parameters could be cumbersome. As introduced in former subsections, there are 16 learnable parameters designed in the WB-Model that control the hydrological behaviors. The ranges of these parameters need to be adjusted manually until the hybrid model produces the most accurate results. For example, the ranges of hydrometeorological parameters that introduced from former models (Jiang et al., 2020; Patil and Stieglitz, 2014), such as T_{min} , T_{max} , f, etc., were set for all catchments while other parameters, especially the geological related and storage related parameters introduced in this study, such as S_{pmax} , Q_{pmax} , S_{cmax} , S_{min} , etc., should be adjust for specific catchments to reach to best simulation performance. Consequently, the parameters calibration process of hybrid model for some catchments could cost lots of time to achieve the best simulation results.

706 The second limitation is that due to the lack of detailed geographic information, the lumped 707 GWL model cannot be applicable to all basins. Supporting by the universality of water balance constrains and the powerful computing ability of DLA, our expectation of this hybrid model is that 708 709 the model can simulate GWL at any point in the basin accurately if the hydrometeorological data 710 is sufficient and the historical GWL observation is provided. This is because we believe the hybrid 711 model can perceive the geological information to give accurate GWL simulations by adjusting learnable parameters set of WB-Model. However, in this study, the hybrid model failed to improve 712 713 the simulation performances for 20 out of 91 catchments. This is because the lack of detailed 714 geographic information makes some groundwater recharge and discharge processes, such as groundwater recharge from surrounding confined aquifers, impossible to be reflected through 715 single water balance constrains. Therefore, the proposed hybrid model cannot be applied to the 716 simulation of GWL in all catchments, especially in those with complex groundwater recharge and 717 discharge process. This is also reflected in the simulation performance of WB-Model: although 718 WB-Model can simulate the fluctuation state of GWL, most of the simulation results are not 719 accurate, and accurate simulation results are obtained by embedding the preliminary result into 720 another DL structure. 721

The third limitation is the black box nature of proposed hybrid model remains unrevealed. The main merit of adding water balance ODEs into the DL architecture lies in improving the awareness of hydrological knowledge of hybrid model by guiding WB-Model to capture the fluctuation patterns of GWL. However, adding physics is not enough to reduce the black box nature of DL model (Barredo Arrieta et al., 2020). Despite we have successfully combined DL
with physical constrains, making the hybrid model present extremely powerful forecasting ability,
we have not found further interpretable physical meaning from the proposed hybrid model.
Especially, for the two-layer 1D-CNN model in the hybrid model, which only provides the
outstand non-linear fitting ability, the fitting process remains uninterpreted from hydrological
perspective in this study.

The fourth limitation is that the two- or three-dimensional processes of groundwater flow are 732 not considered in this study. In general, the study of groundwater simulations is always related 733 734 with two or three-dimensional groundwater flow with specific boundary conditions and governing equations. Moreover, detailed groundwater simulation is usually focusing on find numerical or 735 analytical solutions of two or three-dimensional PDEs related to groundwater by numerical method 736 or analytical method (Liang et al., 2018). The hybrid model we proposed is based on the 737 combination of DLA and vertical one-dimensional water balance equation in catchment scale. 738 Although it is a breakthrough to the traditional research methods, it cannot describe the two and 739 three-dimensional groundwater flow processes, which might make it not delicate or convincing 740 enough from the perspective of traditional groundwater flow simulation community. In other 741 742 words, we embedded physical constraints to improve the simulation performance of DL model, but the physical constraints are not detailed enough compared with traditional methods. More 743 specifically, although water balance constraints with 16 self-learnable parameters are well 744 745 implanted into WB-model, this number is far less than prevalent DL model, such as 1D-CNN model, which in our case, has 568 parameters in first 1D-CNN layer and 9 parameters in the second 746 747 1D-CNN layer. Parsimonious DL model does not ensure a learned model will capture all important 748 information in the data sensed about the external world (Ma et al., 2022). As a result, lacking

complexity is a main reason why WB-Model cannot outperform the GRU model in terms of theaccuracy of GWL simulation.

Based on the limitations mentioned above, we believe that this study has potential for further 751 improvements. Firstly, more detailed groundwater-related processes, such as human activities, 752 could be implanted into WB-Model. The proposed hybrid model could be more flexible if multiple 753 potential processes were wrapped into WB-Model for options. In that case, the hybrid model would 754 decide which processes are suitable to be considered for specific catchments by adjusting the 755 parameters and super parameters of WB-Model. This would improve the adaptability of hybrid 756 757 model greatly if the groundwater-related processes could be considered comprehensively. Secondly, cutting-edge DL interpretive methods, such as expected gradients (EG) (Erion et al., 758 2021), could be applied to decipher the machine-captured patterns and inner workings of the hybrid 759 model. The main purpose of such study will be revealing the black-box process of the proposed 760 hybrid model, so that the hydrological cognition from the perspective of machine learning models 761 under physical constraints could be obtained to facilitate our improving our understanding for 762 specific hydrological processes. Thirdly, the threshold mechanism of hydrological models in our 763 WB-Model framework might be further improved by the idea of gating mechanism from GRU or 764 765 LSTM model. Lastly, to optimize the physical constrains from one-dimensional ODEs to higher dimensional PDEs. The difficulty of groundwater simulation is much higher than that of surface 766 water simulation because of its complex two and three-dimensional flow characteristics (Wang et 767 768 al., 2021). Limited by the lack of geophysical data, the lumped hybrid model based on water balance mechanism has not maximized the power of integrating DLAs and physical constrains. 769 Therefore, the combination of PDEs-governed distributed groundwater flow processes and DL 770 771 model will be an important research direction in the future. Regardless, the proposed hybrid model is already the most novel model in terms of predicting GWL with time series hydrometeorologicalinputs.

774 4 Conclusions

In this study, a lumped hybrid groundwater model with water balance equations as physical constraint and DL methods as core algorithm is proposed to simulate the fluctuation of GWL. In the hybrid model, water balance guided OEDs were wrapped into a WB-Model by self-designing the specific algorithms of RNN model. We tested the model with the CAMELS dataset from 91 catchments located in the middle eastern CONUS, and two pure DL models, 1D-CNN and GRU, were established for comparison of simulation performance. The main findings of this study are summarized as follows:

First, the hybrid model presented high accuracy of simulating the fluctuation of GWL without 782 using detailed hydrogeological information of the catchments. The preliminary GWL simulated 783 from the WB-Model enhances the learning ability of hybrid model. Consequently, the physics 784 constrained DL model outperformed the pure DL models significantly in 71 out of 91 catchments 785 in this study. Moreover, comparing with traditional distributed GWL simulation models, DLA 786 reduces the cost of data as well as the difficulty of model setup while still provides accurate 787 788 simulation results in perspective of traditional standards, which makes the hybrid model more suitable for GWL predictions in less gauged or ungauged basins. 789

Second, the self-designed RNN model with water balance constraints proposed in this study embeds the main groundwater-related water balance formulas, which are referred to the traditional distributed hydrological models, into the recurrent neutral networks. The specific values of the parameters in water balance related formula and water storages in designed reservoirs are determined by the hybrid model through iterative algorithms, activation functions and loss functions of DLAs. Consequently, this equips our hybrid model with the ability to learn the groundwater-related water allocation processes at the catchment scale. In conclusion, the physical constrained hybrid model presents better adaptability and generalization ability comparing with the pure DL models.

Last, compared with the pure DL model, the hybrid DL model proposed has better robustness. This is reflected from the fact that the hybrid models outperformed pure DL model with different strategies of input features. More specifically, the NSE values of simulations from hybrid model have higher values and lower fluctuation range. The reason is that the preliminary simulation results provided by the physically constrained WB-Model strengthen the learning ability of DLA for groundwater fluctuation characteristics, so that the hybrid model can be affected by potential noise data as little as possible.

This study shows the superiority and powerful simulation ability of DL model based on 806 physical constraints. We have abandoned the general idea of building traditional groundwater 807 models: solving two-dimensional or three-dimensional groundwater flow problems through the 808 iterative method of distributed models or finding analytical solutions. Instead, we embedded water 809 balance equations into DLA for regional GWL simulation. As a result, this hybrid model presented 810 811 great accuracy, adaptability, generalization ability, and robustness even without detailed geological data of catchments, which demonstrated the possibility of application of proposed 812 model for GWL simulation in ungauged or lack of gauged catchments. Although there are many 813 814 limitations and potential improvements for the proposed model, we believe that the general performance of the proposed model would increase trust in data-driven approaches on hydrological 815 816 modellings especially when physical constraints related to hydrological sciences are integrated 817 with DLAs.

818	Acknowledgment	S
-----	----------------	---

- This study was supported by the National Natural Science Foundation of China (51909117)
- and Natural Science Foundation of Shenzhen (JCYJ20210324105014039).

821 Author contributions

- 822 Conceptualization: Haiyun Shi
- 823 Methodology: Hejiang Cai, Haiyun Shi
- Formal analysis: Hejiang Cai
- 825 Validation: Hejiang Cai, Zhaoqiang Zhou, Shijie Jiang
- 826 Software: Hejiang Cai, Zhaoqiang Zhou
- 827 Supervision: Haiyun Shi, Suning Liu, Vladan Babovic
- 828 Writing original draft: Hejiang Cai
- 829 Writing review & editing: Haiyun Shi, Suning Liu, Vladan Babovic, Shijie Jiang
- Funding Acquisition: Haiyun Shi

831 Data availability and coding statement

- The groundwater level data are openly available at <u>https://waterdata.usgs.gov/nwis/gw</u>, and
- the CAMELS dataset is freely available at <u>https://ncar.ucar.edu</u>.
- We gratefully acknowledge the use of the dataset provided by Addor et al. (2017) and
- Newman et al. (2015) and the codes of P-RNN model publicly provided by Jiang et al. (2020).

836 **References**

- Addor, N., A. J. Newman, N. Mizukami, and M. P. Clark (2017), The CAMELS data set: catchment attributes and
 meteorology for large-sample studies, *Hydrology and Earth System Sciences*, *21*(10), 5293-5313.
- Ayzel, G., Heistermann, M., 2021. The effect of calibration data length on the performance of a conceptual
- 840 hydrological model versus LSTM and GRU: a case study for six basins from the CAMELS dataset. Comput. Geosci.
- 841 149 https://doi.org/10.1016/j.cageo.2021.104708.

- Babovic, V. (2009), Introducing knowledge into learning based on genetic programming, *Journal of Hydroinformatics*, *11*(3-4), 181-193.
- Babovic, V., and M. Keijzer (2002), Rainfall runoff modelling based on genetic programming, *Nordic Hydrology*,
 33(5), 331-346.
- 846 Banwart, S. A., N. P. Nikolaidis, Y.-G. Zhu, C. L. Peacock, and D. L. Sparks (2019), Soil Functions: Connecting
- Earth's Critical Zone, Annual Review of Earth and Planetary Sciences, 47(1), 333-359.
- Barredo Arrieta, A. et al., 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and
 challenges toward responsible AI. Information Fusion, 58: 82-115.
- 850 Beven, K. J., and H. L. Cloke (2012), Comment on "Hyperresolution global land surface modeling: Meeting a grand
- challenge for monitoring Earth's terrestrial water" by Eric F. Wood et al, *Water Resources Research*, 48(1).
- 852 Broxton, P. D., W. J. D. Leeuwen, and J. A. Biederman (2019), Improving Snow Water Equivalent Maps With Machine
- Learning of Snow Survey and Lidar Measurements, *Water Resources Research*, 55(5), 3739-3757.
- Butcher, J. C. (2000). Numerical methods for ordinary differential equations in the 20th century. Journal of Computational and Applied Mathematics, 125(1-2), 1-29.
- 856 Cai, H., H. Shi, S. Liu, and V. Babovic (2021), Impacts of regional characteristics on improving the accuracy of
- 857 groundwater level prediction using machine learning: The case of central eastern continental United States, *Journal*
- 858 of Hydrology: Regional Studies, 37.
- Chadalawada, J., V. Havlicek, and V. Babovic (2017), A Genetic Programming Approach to System Identification of
 Rainfall-Runoff Models, *Water Resources Management*, *31*(12), 3975-3992.
- 861 Chadalawada, J., H. M. V. V. Herath, and V. Babovic (2020), Hydrologically Informed Machine Learning for
- Rainfall Runoff Modeling: A Genetic Programming Based Toolkit for Automatic Model Induction, *Water Resources Research*, 56(4).
- Chen, C., W. He, H. Zhou, Y. Xue, and M. Zhu (2020), A comparative study among machine learning and numerical
- models for simulating groundwater dynamics in the Heihe River Basin, northwestern China, Sci Rep, 10(1), 3904.
- 866 Clark, M. P., et al. (2015a), Improving the representation of hydrologic processes in Earth System Models, *Water*
- 867 *Resources Research*, 51(8), 5929-5956.
- Clark, M. P., et al. (2015b), A unified approach for process Dased hydrologic modeling: 2. Model implementation
- and case studies, *Water Resources Research*, 51(4), 2515-2542.

- 870 Cuthbert, M. O., T. Gleeson, N. Moosdorf, K. M. Befus, A. Schneider, J. Hartmann, and B. Lehner (2019), Global
- patterns and dynamics of climate–groundwater interactions, *Nature Climate Change*, 9(2), 137-141.
- 872 Dunkerley, D. L., and T. L. Booth (1999), Plant canopy interception of rainfall and its significance in a banded
- landscape, arid western New South Wales, Australia, *Water Resources Research*, 35(5), 1581-1586.
- 874 Erion, G., Janizek, J. D., Sturmfels, P., Lundberg, S. M., & Lee, S. I. (2021). Improving performance of deep learning
- models with axiomatic attribution priors and expected gradients. Nature Machine Intelligence, 3(7), 620–631.
- Fan, Y., H. Li, and G. Miguez-Macho (2013), Global patterns of groundwater table depth, *Science*, *339*(6122), 940943.
- Fatichi, S., D. Or, R. Walko, H. Vereecken, M. H. Young, T. A. Ghezzehei, T. Hengl, S. Kollet, N. Agam, and R.
- Avissar (2020), Soil structure is an important omission in Earth System Models, *Nat Commun*, *11*(1), 522.
- 880 Fatichi, S., et al. (2016), An overview of current applications, challenges, and future trends in distributed process-
- based models in hydrology, *Journal of Hydrology*, 537, 45-60.
- Feng, D., K. Fang, and C. Shen (2020), Enhancing Streamflow Forecast and Extracting Insights Using Long □ Short
- Term Memory Networks With Data Integration at Continental Scales, *Water Resources Research*, 56(9).
- Frame, J, Ulrich, P, Nearing, G, Gupta, H, Kratzert, F, (2022). On strictly enforced mass conservation constraints for
- 885 modeling the rainfall-runoff process, EarthArXiv, doi: <u>https://doi.org/10.31223/X5BH0P</u>
- 886 Gao, S., Y. Huang, S. Zhang, J. Han, G. Wang, M. Zhang, and Q. Lin (2020), Short-term runoff prediction with GRU
- and LSTM networks without requiring time step optimization during sample generation, *Journal of Hydrology*, 589.
- 888 Herrmann, F., N. Baghdadi, M. Blaschek, R. Deidda, R. Duttmann, I. La Jeunesse, H. Sellami, H. Vereecken, and F.
- 889 Wendland (2016), Simulation of future groundwater recharge using a climate model ensemble and SAR-image based
- soil parameter distributions A case study in an intensively-used Mediterranean catchment, Sci Total Environ, 543(Pt
- B), 889-905.
- Hoedt, P.-J., F. Kratzert, D. Klotz, C. Halmich, M. Holzleitner, G. S. Nearing, S. Hochreiter, and G. Klambauer (2021),
- 893 MC-LSTM: Mass-Conserving LSTM, in Proceedings of the 38th International Conference on Machine Learning,
- edited by M. Marina and Z. Tong, pp. 4275--4286, PMLR, Proceedings of Machine Learning Research.
- Horton, R. E. (1933), The Rôle of infiltration in the hydrologic cycle, Eos Trans. AGU, 14(1), 446–460,
- doi:10.1029/TR014i001p00446.
- Jiang, S., Y. Zheng, and D. Solomatine (2020), Improving AI System Awareness of Geoscience Knowledge: Symbiotic

- 898 Integration of Physical Approaches and Deep Learning, *Geophys Res Lett*, 47(13).
- Jiang, S., Y. Zheng, C. Wang, and V. Babovic (2022), Uncovering Flooding Mechanisms Across the Contiguous United
- States Through Interpretive Deep Learning on Representative Catchments, *Water Resources Research*, 58(1),
 e2021WR030185.
- Jiang, S., Y. Zheng, V. Babovic, Y. Tian, and F. Han (2018), A computer vision-based approach to fusing spatiotemporal
 data for hydrological modeling, *Journal of Hydrology*, *567*, 25-40.
- Jiang, S. J., V. Babovic, Y. Zheng, and J. Z. Xiong (2019), Advancing Opportunistic Sensing in Hydrology: A Novel
- Approach to Measuring Rainfall With Ordinary Surveillance Cameras, *WATER RESOURCES RESEARCH*, 55(4),
 3004-3027.
- Jiang, Z.P., H.Y. Shi, S.N. Liu, Z.Q. Zhou, Y. Wang, and H.J. Cai (2021), Evolution characteristics of potential
- evapotranspiration over the Three-River Headwaters Region, *Hydrological Sciences Journal*, *66*(10), 1552-1566.
- 909 Kratzert, F., D. Klotz, M. Herrnegger, A. K. Sampson, S. Hochreiter, and G. S. Nearing (2019), Toward Improved
- 910 Predictions in Ungauged Basins: Exploiting the Power of Machine Learning, Water Resources Research, 55(12),
- 911 11344-11354.
- 912 Kirkby, M.: Hydrograph modelling strategies, in: Processes in Human and Physical Geography, edited by: Peel, R.,
- 913 Chisholm, M., and Haggett, P., Heinemann, London, 69–90, 1975.
- 914 Leavesley, G.H., Lichty, R.W., Troutman, B.M., Saindon, L.G., 1983. Precipitation-Runoff-Modeling-System,
- 915 User's Manual, Water Resource Investigations Report 83-4238, US Geological Survey
- 916 LeCun Y, Boser B, Denker J S, et al. Backpropagation applied to handwritten zip code recognition[J]. Neural
- 917 computation, 1989, 1(4): 541-551.
- 918 LeCun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition[J]. Proceedings of the
- 919 IEEE, 1998, 86(11): 2278-2324.
- Li, L., C. Bao, P. L. Sullivan, S. Brantley, Y. Shi, and C. Duffy (2017), Understanding watershed hydrogeochemistry:
- 921 2. Synchronized hydrological and geochemical processes drive stream chemostatic behavior, *Water Resources* 922 *Research*, 53(3), 2346-2367.
- Liang, X., H. Zhan, and Y.-K. Zhang (2018), Aquifer Recharge Using a Vadose Zone Infiltration Well, *Water Resources Research*, 54(11), 8847-8863.
- 25 Liu, S.N., and H.Y. Shi (2019), A recursive approach to long-term prediction of monthly precipitation using genetic

- programming, *Water Resources Management*, 33(3), 1103-1121.
- Ma, Y., et al. (2022). On the Principles of Parsimony and Self-Consistency for the Emergence of Intelligence.
 https://doi.org/10.48550/arXiv.2207.04630.
- 929 Ma, Y., C. Montzka, B. Bayat, and S. Kollet (2021a), Using Long Short-Term Memory networks to connect water
- table depth anomalies to precipitation anomalies over Europe, Hydrology and Earth System Sciences, 25(6), 3555-
- 931 3575.
- Ma, Y., C. Montzka, B. Bayat, and S. Kollet (2021b), An Indirect Approach Based on Long Short-Term Memory
- 933 Networks to Estimate Groundwater Table Depth Anomalies Across Europe With an Application for Drought
- Analysis, *Frontiers in Water*, *3*.
- 935 Markstrom, S.L., Niswonger, R.G., Regan, R.S., Prudic, D.E., and Barlow, P.M., 2008, GSFLOW-Coupled Ground-
- 936 water and Surface-water FLOW model based on the integration of the Precipitation-Runoff Modeling System
- 937 (PRMS) and the Modular Ground-Water Flow Model (MODFLOW-2005): U.S. Geological Survey Techniques and

938 Methods 6-D1, 240 p.

- Mohan, C., A. W. Western, Y. Wei, and M. Saft (2018), Predicting groundwater recharge for varying land cover and
 climate conditions a global meta-study, *Hydrology and Earth System Sciences*, *22*(5), 2689-2703.
- 941 Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models. Part I-A discussion of
- 942 principles. Journal of Hydrology, 10(3), 282–290.
- 943 Newman, A. J., et al. (2015), Development of a large-sample watershed-scale hydrometeorological data set for the
- 944 contiguous USA: data set characteristics and assessment of regional variability in hydrologic model performance,
- 945 *Hydrology and Earth System Sciences*, 19(1), 209-223.
- Niu, M., Horesh, L., & Chuang, I. (2019). Recurrent neural networks in the eye of differential equations. arXiv e-
- prints, arXiv:1904.12933. Retrieved from <u>https://arxiv.org/abs/1904.12933</u>
- Ossai, I. C., A. Ahmed, A. Hassan, and F. S. Hamid (2020), Remediation of soil and water contaminated with petroleum hydrocarbon: A review, *ENVIRONMENTAL TECHNOLOGY & INNOVATION*, *17*.
- Patil, S., and M. Stieglitz (2014), Modelling daily streamflow at ungauged catchments: what information is necessary?,
- 951 *Hydrological Processes*, 28(3), 1159-1169.
- Raissi, M., Perdikaris, P., Karniadakis, G.E., 2019. Physics-informed neural networks: A deep learning framework for
- solving forward and inverse problems involving nonlinear partial differential equations. J. Comput. Phys. 378, 686-

954 707.

- Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, N. Carvalhais, and Prabhat (2019), Deep learning
 and process understanding for data-driven Earth system science, *Nature*, *566*(7743), 195-204.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back □propagating errors.
 Nature, 323(6088), 533–536.
- 959 R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuglu, P. Kuksa. 2011. Natural Language Processing
- 960 (Almost) from Scratch. Journal of Machine Learning Research 12:2493–2537.
- 961 Sadeghfam, S., A. Ehsanitabar, R. Khatibi, and R. Daneshfaraz (2018), Investigating 'risk' of groundwater drought
- 962 occurrences by using reliability analysis, *Ecological Indicators*, 94, 170-184.
- Sahoo, S., T. A. Russo, J. Elliott, and I. Foster (2017), Machine learning algorithms for modeling groundwater level
 changes in agricultural regions of the U.S, *Water Resources Research*, *53*(5), 3878-3895.
- 965 Su, D., H. Zhang, H. Chen, J. Yi, P.-Y. Chen, and Y. Gao (2018), Is Robustness the Cost of Accuracy? A
- Comprehensive Study on the Robustness of 18 Deep Image Classification Models, Springer International Publishing,Cham.
- Tan, L. S., Z. Zainuddin, and P. Ong (2018), Solving ordinary differential equations using neural networks, edited.
- 769 Tao, W., Q. Wang, L. Guo, H. Lin, X. Chen, Y. Sun, and S. Ning (2020), An enhanced rainfall-runoff model with
- coupled canopy interception, *Hydrological Processes*, *34*(8), 1837-1853.
- 971 Tian, Y., Y. Zheng, B. Wu, X. Wu, J. Liu, and C. Zheng (2015), Modeling surface water-groundwater interaction in
- arid and semi-arid regions with intensive agriculture, Environ Modell Softw, 63, 170-184. Tran, H., E. Leonarduzzi,
- 273 L. De la Fuente, R. B. Hull, V. Bansal, C. Chennault, P. Gentine, P. Melchior, L. E. Condon, and R. M. Maxwell
- 974 (2021), Development of a Deep Learning Emulator for a Distributed Groundwater–Surface Water Model: ParFlow-
- 975 ML, Water, 13(23).
- 976 Trinh, D. H., and T. F. M. Chui (2013), An empirical method for approximating canopy throughfall, *Hydrological*977 *Processes*, 27(12), 1764-1772.
- Viana, F. A. C., R. G. Nascimento, A. Dourado, and Y. A. Yucesan (2021), Estimating model inadequacy in ordinary
 differential equations with physics-informed neural networks, Computers & Structures, 245.
- 980 Wang, N., Zhang, D., Chang, H., & Li, H. (2020). Deep learning of subsurface flow via theory-guided neural network.
- 981 Journal of Hydrology, 584, 124700.

- 982 Wang, N., H. Chang, and D. Zhang (2021), Deep Learning Based Inverse Modeling Approaches: A Subsurface
- 983 Flow Example, Journal of Geophysical Research: Solid Earth, 126(2).
- 984 Yoon, H., S.-C. Jun, Y. Hyun, G.-O. Bae, and K.-K. Lee (2011), A comparative study of artificial neural networks and
- support vector machines for predicting groundwater levels in a coastal aquifer, Journal of Hydrology, 396(1-2), 128138.
- 987 Zhang, J., X. Chen, A. Khan, Y.-k. Zhang, X. Kuang, X. Liang, M. L. Taccari, and J. Nuttall (2021), Daily runoff
- forecasting by deep recursive neural network, *Journal of Hydrology*, 596.
- 989 Zhao, W. L., P. Gentine, M. Reichstein, Y. Zhang, S. Zhou, Y. Q. Wen, C. J. Lin, X. Li, and G. Y. Qiu (2019), Physics-
- 990 Constrained Machine Learning of Evapotranspiration, *Geophys Res Lett*, 46(24), 14496-14507.
- 291 Zhou, Z., H. Shi, Q. Fu, T. Li, T. Y. Gan, and S. Liu (2020), Assessing spatiotemporal characteristics of drought and
- 992 its effects on climate-induced yield of maize in Northeast China, Journal of Hydrology, 588.
- 293 Zipper, S. C., T. Dallemagne, T. Gleeson, T. C. Boerman, and A. Hartmann (2018), Groundwater Pumping Impacts on
- Real Stream Networks: Testing the Performance of Simple Management Tools, *Water Resources Research*, 54(8),
 5471-5486.

1 Author contributions

2	Conceptua	lization:	Haiyun	Shi
-	00110 p t t t			~

- 3 Methodology: Hejiang Cai, Haiyun Shi
- 4 Formal analysis: Hejiang Cai
- 5 Validation: Hejiang Cai, Zhaoqiang Zhou, Shijie Jiang
- 6 Software: Hejiang Cai, Zhaoqiang Zhou
- 7 Supervision: Haiyun Shi, Suning Liu, Vladan Babovic
- 8 Writing original draft: Hejiang Cai
- 9 Writing review & editing: Haiyun Shi, Suning Liu, Vladan Babovic, Shijie Jiang
- 10 Funding Acquisition: Haiyun Shi