# Spatio-temporal machine learning for continental scale terrestrial hydrology

Andrew Bennett<sup>1</sup>, Hoang Tran<sup>2</sup>, Luis De la Fuente<sup>1</sup>, Amanda Triplett<sup>1</sup>, Yueling Ma<sup>3</sup>, Peter Melchior<sup>3</sup>, Reed M. Maxwell<sup>3</sup>, and Laura Elizabeth Condon<sup>1</sup>

<sup>1</sup>University of Arizona <sup>2</sup>Pacific Northwest National Laboratory <sup>3</sup>Princeton University

November 24, 2023

#### Abstract

Integrated hydrologic models can simulate coupled surface and subsurface processes but are computationally expensive to run at high resolutions over large domains. Here we develop a novel deep learning model to emulate continental-scale subsurface flows simulated by the integrated ParFlow-CLM model. We compare convolutional neural networks like ResNet and UNet run autoregressively against our novel architecture called the Forced SpatioTemporal RNN (FSTR). The FSTR model incorporates separate encoding of initial conditions, static parameters, and meteorological forcings, which are fused in a recurrent loop to produce spatiotemporal predictions of groundwater. We evaluate the model architectures on their ability to reproduce 4D pressure heads, water table depths, and surface soil moisture over the contiguous US at 1km resolution and daily time steps over the course of a full water year. The FSTR model shows superior performance to the baseline models, producing stable simulations that capture both seasonal and event-scale dynamics across a wide array of hydroclimatic regimes. The emulators provide over 1000x speedup compared to the original physical model, which will enable new capabilities like uncertainty quantification and data assimilation for integrated hydrologic modeling that were not previously possible. Our results demonstrate the promise of using specialized deep learning architectures like FSTR for emulating complex process-based models without sacrificing fidelity.

#### Hosted file

977690\_0\_art\_file\_11525486\_s33hq9.docx available at https://authorea.com/users/677920/ articles/674778-spatio-temporal-machine-learning-for-continental-scale-terrestrialhydrology

1	
2	Spatio-temporal machine learning for continental scale terrestrial hydrology
3 4	Enter authors here: Andrew Bennett <sup>1</sup> , Hoang Tran <sup>2</sup> , Luis De la Fuente <sup>1</sup> , Amanda Triplett <sup>1</sup> , Yueling Ma <sup>3,4</sup> , Peter Melchior <sup>5,6</sup> , Reed M. Maxwell <sup>3,4</sup> , Laura E. Condon <sup>1</sup>
5	<sup>1</sup> Department of Hydrology and Atmospheric Sciences, University of Arizona, Tucson, AZ, USA.
6 7	<sup>2</sup> Atmospheric Science & Global Change Division, Pacific Northwest National Laboratory, Richland, WA, USA.
8	<sup>3</sup> Department of Civil and Environmental Engineering, Princeton University, Princeton, NJ, USA
9	<sup>4</sup> High Meadows Environmental Institute, Princeton University, Princeton, NJ, USA.
10	<sup>5</sup> Department of Astrophysical Sciences, Princeton University, Princeton, NJ, USA.
11	<sup>6</sup> Center for Statistics and Machine Learning, Princeton University, Princeton, NJ, USA.
12	
13	Corresponding author: Andrew Bennett (andrbenn@arizona.edu)
14	Key Points:
15 16	• Deep learning models can emulate a complex integrated hydrology model that simulates groundwater over the United States
17 18	• We developed a new model architecture that is more robust over longer simulation periods than off-the-shelf neural networks
19 20 21	• Deep-learning based emulators of complex models enable new applications such as real- time forecasting and estimating uncertainties

## 22 Abstract

23 Integrated hydrologic models can simulate coupled surface and subsurface processes but are

computationally expensive to run at high resolutions over large domains. Here we develop a novel deep

25 learning model to emulate continental-scale subsurface flows simulated by the integrated ParFlow-CLM

26 model. We compare convolutional neural networks like ResNet and UNet run autoregressively against

our novel architecture called the Forced SpatioTemporal RNN (FSTR). The FSTR model incorporates

separate encoding of initial conditions, static parameters, and meteorological forcings, which are fused in

a recurrent loop to produce spatiotemporal predictions of groundwater. We evaluate the model architectures on their ability to reproduce 4D pressure heads, water table depths, and surface soil moisture

over the contiguous US at 1km resolution and daily time steps over the course of a full water year. The

FSTR model shows superior performance to the baseline models, producing stable simulations that

capture both seasonal and event-scale dynamics across a wide array of hydroclimatic regimes. The

34 emulators provide over 1000x speedup compared to the original physical model, which will enable new

35 capabilities like uncertainty quantification and data assimilation for integrated hydrologic modeling that 36 were not previously possible. Our results demonstrate the promise of using specialized deep learning

36 were not previously possible. Our results demonstrate the promise of using specialized deep learning 37 architectures like FSTR for emulating complex process-based models without sacrificing fidelity.

38

# 39 Plain Language Summary

40 Computational models are important for understanding and predicting terrestrial hydrology, but

41 our most physically detailed models can be time-consuming and expensive to run over large

regions. In this study, we trained deep learning models to emulate a complex hydrology model

that simulates groundwater flow over the contiguous US. We developed a new model

44 architecture called FSTR that captures spatiotemporal patterns better than standard deep learning

45 models. FSTR is over 1000 times faster than ParFlow, the original hydrologic model. This

enables new possibilities like forecasting groundwater changes and estimating uncertainties. Our

47 results show that specialized deep learning architectures can accurately emulate complex

48 hydrologic models while drastically reducing computation time.

# 49 **1 Introduction**

50 Computational models have been hugely successful in predicting and building understand of

51 Earth and environmental systems. Since their early use in the mid twentieth century these models

52 have increased in complexity to account for higher spatiotemporal resolution and physical

53 process complexity. To achieve the highest possible degree of physical realism scientists and

<sup>54</sup> engineers often must run these models on high-performance computing clusters and

supercomputers, which require large institutional investment to build and maintain. In addition,

the computational complexity of using these models on such platforms requires considerable

57 expertise to use effectively. As a result, these large-scale and highly complex models are

58 typically only used by a small subset of researchers.

59 Emulation (also referred to as reduced order models (ROMs) or surrogate models) has long

been a popular method to reduce the computational complexity of running simulations in a

number of domains (Astrid et al., 2008; Razavi et al., 2012; C. Wang et al., 2014). Reducing the

62 computational complexity of process based models makes it possible to build more complex

63 workflows such as building model chains, performing parameter calibration or sensitivity

experiments (Cheng et al., 2023), and running more scenarios to better understand uncertainties

65 (Kasim et al., 2021). One avenue that is becoming increasingly popular is the use of deep

learning (DL) based methods for building emulators of process based models (Doury et al., 2023;
Leonarduzzi et al., 2022; Reichstein et al., 2019; Tran et al., 2021).

Deep learning has recently and quickly become a standard piece of the computational

modeling toolkit and offers almost universal applicability to modeling tasks (Jordan & Mitchell,

2015). It has also proven very powerful in allowing researchers to take models from disparate

applications and apply them to new problems (Khan et al., 2022). This is true in the Earth system

sciences and related fields, where deep learning models have been used to simulate atmospheric

phenomena (Brenowitz et al., 2020), predict flow in rivers (Kratzert et al., 2018), and monitor

<sup>74</sup> land use (Xu et al., 2017) among many other applications. Specifically in hydrology the majority

of applications are based around streamflow modeling or other forms of "point-scale"

applications where the models are trained on individual sites (Bennett & Nijssen, 2021; de la

Fuente et al., 2023; Gauch et al., 2021). However, we know that in subsurface hydrology lateral

flow occurs, and can have large impacts on both human and natural systems (Condon &
Maxwell, 2019; Fan, 2015).

80 Hydrologic models that treat both the coupled surface and groundwater systems as well as 81 account for lateral flow in the subsurface are commonly referred to as integrated hydrologic

models. These models are among the most complex and comprehensive representations of the

terrestrial hydrologic cycle that have been developed. However, they are often difficult to run

84 because they are data-hungry and computationally heavy. This is particularly true in the

subsurface, where observations are sparse and parameters are hard to measure (Blöschl et al.,

86 2019). Even when data is available, it is often difficult to calibrate these models because of the

87 computational complexity and large-dimensional search space over parameter configurations

88 (O'Neill et al., 2021). It should also be noted that the lack of spatiotemporally complete

observations/reanalysis data for groundwater is a large reason that purely data-driven approaches
 have not emerged as they have in the weather forecasting domain (Chen et al., 2023; Keisler,

have not emerged as they have in the weather forecasting domain (Chen et al., 2023; Keisler,
2022; Lam et al., 2022). Because of these challenges the use of emulator or surrogate models is

an appealing approach to improving the usability of integrated hydrologic models.

In this study we demonstrate the use of modern deep learning based emulators of a continental scale integrated hydrologic model, without sacrificing spatiotemporal or process fidelity. Specifically, we emulate subsurface flow of the ParFlow-CLM model, developed over a large portion of the contiguous US at a high spatiotemporal resolution (Maxwell et al., 2015; Maxwell & Condon, 2016; O'Neill et al., 2021). Previous work has shown that deep learning is an effective approach to this problem, showing good performance on synthetic benchmarks ( Maxwell et al., 2021)and on smaller domains (Leonarduzzi et al., 2022; Tran et al., 2021).

In this work, we compare the ability of three different deep learning model architectures to emulate 3d subsurface pressure fields simulated by ParFlow. We compare ResNet, UNet, and

Forced-SpatioTemporal-RNN (FSTR) architectures; the first two are off-the-shelf convolutional neural networks (CNN) that we apply in an autoregressive fashion to build up spatiotemporal predictions. That is, we feed the previous output of the model as an input to the next step of the prediction process, in addition to other features.

The FSTR model is a novel adaptation of the PredRNN model with action-conditioning, which is a video-prediction model that has proven capable for atmospheric & hydrologic modeling as well as robotics (Tran et al., 2021; Wang et al., 2017). We make use of the robotics terminology of "action-conditioning" to explicitly account for the meteorological forcings acting on the hydrology in a way that is separate from the subsurface parameters/geology and the initial

111 state of the domain that we are simulating.

We compare the performance of the three model architectures for emulating the evolution of

- both pressure heads at multiple depths as well as the derived soil moisture states and water table
- depths. We find that the ResNet produces unstable results over long simulation rollouts, while
- the UNet and FSTR both show good overall capabilities at matching spatial and temporal
- patterns. Our FSTR architecture consistently shows the best performance results, and is capable
- of simulating a year of the entire domain in less than an hour on a single 40 GB Nvidia A100
   GPU, showing a >1000 times speedup as compared to the original simulations run on >3000
- 119 CPU cores. Based on these findings we believe that our FSTR architecture could form the basis
- for a new set of modeling capabilities to fine tune model parameters and enable real-time
- 121 ensemble-based forecasting.

# 122 **2 Methods**

123

2.1 Modeling domain and data

Our modeling domain is based on the ParFlow CONUS1.0 model which is documented by

125 (O'Neill et al., 2021). The domain covers the majority of the contiguous United States, plus

some small portions of Canada and Mexico (Figure 1). We are primarily interested in

representing the subsurface hydrology in this domain at a 1km gridded spatial resolution, with a

daily timestep. The gridded domain amounts to a regular grid of 3342 by 1888 km in the

129 longitudinal and latitudinal directions, respectively. The depth layers of our simulations increase

in the downward direction, starting with shallow surface layers and a large groundwater layer.

The depths of each layer are 0.1, 0.3, 0.6, 1, and 100 m from the surface to the bottom for a total of 5 layers.

The simulations that we use for training/validation/testing have previously been validated against many observational datasets across multiple variables (e.g. streamflow,

evapotranspiration, snow) with favorable results given the default parameter sets (O'Neill et al.,

136 2021). These simulations cover water years 2003-2006 at an hourly timescale. We aggregate all

the data to a daily timescale by taking the daily means, totals, minimums, and maximums where appropriate.

139



140

Figure 1. The extent of our modeling domain, covering most of the Contiguous United States
 (CONUS). We show the long-term average water table depth to highlight the spatial variability
 of the domain.

In this study we are primarily concerned with modeling subsurface and surface water flow, which is represented in ParFlow with Richard's equation (Richards, 1931) parameterized by the van Genuchten closure relations between pressure heads and saturation content (van Genuchten, 1980). Specifically, our emulation target is the full four-dimensional (time+space) pressure head field. From this pressure head field we can use the closure equations and additional calculations to calculate soil moisture and water table depth. The governing equations that describe the dynamics of the subsurface flows are given as:

$$S_s S(\psi) \frac{\partial \psi}{\partial t} + \phi \frac{\partial S(\psi)}{\partial t} = \Delta \left( -K_s(x)k(\psi) \cdot \nabla(\psi - z) \right) + q$$

152

153 Where  $S_s$  is the specific storage  $[L^{-1}]$ , S is the relative saturation [-],  $\psi$  is the pressure head 154 [L],  $K_s$  is the saturated hydraulic conductivity  $[LT^{-1}]$ , k is the relative permeability [-],  $\phi$  is the 155 porosity [-], and q is a source-sink term $[L^3T^{-1}]$ .

# 156 2.2 Model architectures

157 We explore three deep learning model architectures to emulate the simulations described 158 previously. Here we will describe the overall structure of each of these models, but will leave the

detailed input/output setup of them for the following section. We employ three model

architectures in this study: a Residual Network (ResNet; He et al., 2015), a UNet (Ronneberger

161 et al., 2015), and a newly developed architecture that we refer to as the

162 ForcedSpatioTemporalRNN (FSTR) model. The first two of these are relatively "standard"

models in the deep learning literature at this point and have been used for a large array of tasks.

164 Their architectures of the ResNet and UNet are shown in schematic form in figure 2, while the

165 FSTR architecture is shown in figure 3.

The development of the ResNet architecture is considered a milestone in the development of 166 machine learning to the task of image classification and paved the way for the modern deep 167 learning revolution (He et al., 2015). We consider the ResNet a strong baseline architecture to 168 compare against as has been adopted by other studies with similar approaches (Haber & 169 Ruthotto, 2018; Kochkov et al., 2021; Ott et al., 2020). This model consists of stacks of "residual 170 blocks", usually convolutional layers followed by a nonlinear activation function. Following 171 these stacks is a residual connection, which adds the input back to the output of the residual 172 block. Architectures with residual connections have been found to train more effectively, 173 especially in very deep networks. We stack two residual blocks, each consisting of 1 depthwise-174 separable convolutional layer with a layer norm and activation function, and a final 175

convolutional layer. We use four layers with each layer having a hidden dimension of 256channels for this study.

The second architecture that we consider is the UNet, which is named as such because of its 178 use of successive downsampling and upsampling layers (along with skip connections) which 179 allow the model to capture spatial relationships at varying resolutions. This type of architecture 180 is considered state of the art in image segmentation which has applications in both remote 181 sensing (Yuan et al., 2021) and medical imaging (Ronneberger et al., 2015). Like the ResNet, the 182 UNet is mainly composed of stacks of convolutional layers. However, the UNet makes use of 183 downsampling and upsampling to get a "multi-resolution" view of the data. In our model, we use 184 185 stacks of these downsampling and upsampling layers which are composed of convolutional layers that are either preceded by MaxPool layers or followed by bilinear interpolation layers for 186 down and upsampling, respectively. Each downsampling and upsampling layer consists of two 187 convolutions followed by an activation. At parallel levels in the downsampling/upsampling skip 188 connections are used to transfer information at multiple resolutions to the upsampling, which 189 helps to preserve spatial structure in the data. Here we set the base dimension for the 190 convolutional layer to be eight at the input and output and double/halve the hidden dimension for 191

192 each down/up sampling layer respectively.





Figure 2. Diagrams of the two baseline model architectures used in our experimental setup. Left
 shows a ResNet architecture, which stacks convolutional blocks with residual connections that

196 propagate the input signal into deeper layers. The right panel shows the UNet architecture, which 197 also consists of convolutional blocks, but differs from the ResNet by using them in successive

- downsampling and upsampling configurations, which consider multiple resolutions of the input
- 199

We also develop a novel architecture based on the PredRNN model, which is itself based on 200 the Convolutional LSTM model (ConvLSTM; Shi et al., 2015). It attempts to take best practices 201 from image modeling via CNNs and sequence modeling from recurrent neural networks (RNNs), 202 particularly with the application of Long Short Term Memory networks (LSTMs; Hochreiter & 203 Schmidhuber, 1997). The key insights that the developers of PredRNN and associated models 204 had over the ConvLSTM architectures was that an additional memory/hidden state would better 205 reflect the spatiotemporal state of the system, and could be shared amongst model layers. This, 206 along with several other procedural training techniques led the PredRNN model to be one of the 207 most performant video prediction models available (Wang et al., 2017). In the original PredRNN 208 paper, the authors also introduced an "Action-Conditioned" variant, which allows the input 209 sequence of a robotic arm to be used as a part of the prediction algorithm. We make use of this 210 modification because the input to the robotic system "acts" on the video stream in a similar way 211 that meteorological forcings "act" on the evolution of hydrologic states. 212

The main modification that we make to the PredRNN structure is the use of encoders to 213 initialize the memory and cell states for the model. These states are updated throughout the 214 recurrent loop, and by default, are initialized as zeros in the PredRNN structure. However, our 215 insight is that the initial conditions and subsurface parameters can be considered a sort of 216 217 byproduct of the true memory of the natural environment, and thus can be used to initialize these hidden states. We use the initial conditions (i.e. the 3-dimensional pressure heads for the domain 218 being simulated) to initialize the memory state and the parameter values (e.g. porosity and 219 permeability for each cell) to initialize the cell states. Both encoders consist of convolutional 220 layers that project the inputs into a higher dimensional space that matches the hidden states of the 221 Action-Conditioned ST-LSTM which forms the backbone of the FSTR model. The initial 222

conditions are used as input to the memory encoder as well as used as the starting input pressure

field to the model. We use two layers of the AC-ST-LSTM layer in our model, each having a

hidden dimension of 64.



data.

Figure 3. An architecture diagram of our proposed model architecture, the Forced Spatio 227 Temporal RNN (FSTR). Tensor variables are represented in yellow colors, while neural-network 228 layers with trainable weights are drawn in blue. The red "initialization" phase is only run a single 229 time per training example, while the update arrow is run in a recurrent loop. Quantities in the 230 orange outlined boxes represent the hidden states and model inputs. 231

232

The core of the FSTR model is the Action-Conditioned SpatioTemporal LSTM layer (AC-233

ST-LSTM), which was introduced in Wang et al. (2017). This layer modifies the ConvLSTM 234 model in ways that allow it to take in external inputs on the system in the form of an action that 235

is fused to the hidden state by an elementwise multiplication. The equations for the AC-ST-236 237 LSTM layer are given by

$$g_{t} = tanh(W_{xg} * X_{t} + W_{hg} * H_{t-1}^{l})$$

$$i_{t} = \sigma(W_{xi} * X_{t} + W_{hi} * H_{t-1}^{l})$$

$$f_{t} = \sigma(W_{xf} * X_{t} + W_{hf} * H_{t-1}^{l})$$

$$C_{t}^{l} = f_{t} \odot C_{t-1}^{l} + i_{t} \odot g_{t}$$

$$g_{t}' = tanh(W_{xi}' * X_{t} + W_{mg} * M_{t}^{l-1})$$

$$i_{t}' = \sigma(W_{xi}' * X_{t} + W_{mf} * M_{t}^{l-1})$$

$$f_{t}' = \sigma(W_{xf}' * X_{t} + W_{mf} * M_{t}^{l-1})$$

$$M_{t}^{l} = f_{t}' \odot M_{t}^{l-1} + i_{t}' \odot g_{t}'$$

$$o_{t} = \sigma(W_{xo} * X_{t} + W_{ho} * H_{t-1}^{l} + W_{co} * C_{t}^{l} + W_{mo} * M_{t}^{l})$$

$$H_{t}^{l} = o_{t} \odot tanh(W_{1\times1} * [C_{t}^{l}, M_{t}^{l}])$$

$$V_{t}^{l} = (W_{hv} * H_{t-1}^{l}) \odot (W_{av} * A_{t-1})$$

Where  $X_t$  is the input,  $W_{t}$  are the weight matrices,  $H_t^l$  is the hidden state,  $C_t^l$  is the cell state,  $M_t^l$ 238

is the spatiotemporal memory state, A is the action tensor, i, f, g are the gating functions (along 239

with their primed counterparts),  $V_t^l$  is the action fusion which allows for external inputs to 240

modify the system, and  $o_t$  is the output for timestep t, and l is the depth layer of the network. 241

Finally, the update step for the AC-ST-LSTM layer is 242

# $H_t^l, C_t^l, M_t^l = ACSTLSTM(X_t, V_t^l, C_{t-1}^l, M_t^{l-1})$

Our model structure considers the meteorological forcings to be "action" tensors rather than 243 model inputs directly. This delineation of information partitioning between the initial conditions, 244 245 parameter values, and meteorological forcings in FSTR is very similar to how ParFlow and many other hydrologic models operate. 246

All models take in daily minimum temperature, daily maximum temperature, total 247 precipitation, snowmelt, and bulk evapotranspiration as the boundary condition forcing. Static 248 parameters that the models take as input are porosity, permeability, van Genuchten n, van 249 Genuchten alpha, topographic index, elevation, and a measure of distance to the nearest stream 250 based on the digital elevation map. The static parameters were chosen to include the major 251 parameters required to solve Richard's equation as well as represent major topographic features 252 at both point and aggregated scales. The models output 4-dimensional spatiotemporal predictions 253 254 of the subsurface pressure heads for all grid cells in the input domain. All models are run in an auto-regressive fashion, meaning they start with the initial conditions as a starting point and then 255 evolve their output in response to the forcings and parameter value inputs, at which point they 256

use their own prediction as the initial state for the next time step. 257

At inference time (that is, after the models have been trained) we use the models to produce 258 the full 4D pressure head field, which is the quantity that is treated most fundamentally by the 259 form of Richard's equation that is used in the ParFlow simulations. However, when using 260 ParFlow model outputs it is often more useful to calculate other quantities such as soil moisture 261 and water table depth from the subsurface pressure head field. We make a similar conversion to 262 the model outputs by processing them into water table depth and surface soil moisture. These 263 quantities are calculated in the same way that they are in standalone ParFlow simulations, but 264 rather than using the functionality to compute this translation from ParFlow we re-implemented 265 the routines in a PyTorch compatible layer, which in theory could provide these translations at 266 training time. We do not take this approach here, however, due to several technical challenges 267 which will be discussed later but could form the basis for training a fully differentiable model 268 directly to data in the future. 269

### 270 2.3 Experimental setup

In this study we explore the ability of the three neural network architectures to reproduce the 271 3d ParFlow pressure heads across the CONUS domain. To quantify this, we train each model on 272 water years 2003 and 2004 and validate them on water year 2005. The testing dataset that we 273 evaluate against is from the water year 2006, and all results shown here are from this set. As 274 mentioned previously, we aggregate variables to daily timesteps to reduce the overall complexity 275 of the data. While we are interested in building emulators that maintain high spatiotemporal 276 resolution we are primarily focused on seasonal to annual modeling, reflecting timescales that 277 278 groundwater interactions tend to take place at.

We train all models in two phases. For the first phase we use a one-cycle learning rate scheduler with a maximum learning rate of 1e-3. We train on square chips of 64 by 64 pixels in size, and rollout horizons of 35 timesteps. This balance between medium spatial sizes and time horizon provides the model a good baseline for matching both spatial and temporal patterns, which resulted in the best overall training performance. During this portion of the training we augment the L2 loss with an additional term that penalizes gradients in the spatial directions of the predictions (Serifi et al., 2021). This modified loss function is given as:

286

$$\mathcal{L}(y,\hat{y}) = |y - \hat{y}|_2 + |\nabla y - \nabla \hat{y}|_2$$

Following an initial training epoch on this setup we then set the loss function to a pure L2 287 loss and continue training. At this point we train for another epoch using 48 pixel chips and a 288 rollout horizon of 90 days. This helps stabilize the autoregressive loop of the model, as well as 289 learn seasonal trends. Following this, we then train a final epoch using a 14-day rollout and 48 290 pixel chips. This final epoch of training is aimed at improving the model's ability to reproduce 291 fast dynamics, particularly around responding to individual storm events. We found that training 292 across all three phases improved the model prediction in terms of both stability of the predictions 293 and the overall accuracy for all model types. However, the final training stage using the shorter 294 rollouts helped the least. We also found that changing the ordering of the training phases led to 295 degradation in performance. All models were trained on a single Nvidia A100 40GB GPU. 296 Training times varied by model, but generally took roughly 50 hours each, with the ResNet being 297 cheapest and FSTR being the most expensive. 298

## 299 **3 Results**

To analyze the performance of each of the emulator architectures we calculated the overall 300 root mean square error (RMSE) and Pearson correlation (hereafter, just correlation) for each grid 301 cell in the domain with respect to the original simulations. The resulting spatial maps for these 302 measures are shown in figure 4. From figure 4a we see that the ResNet shows pockets of high 303 304 error, particularly in the Western portion of the domain, but also in the upper midwest and southern central regions. These spots of high error are due to model instability at the full lead-305 time. We will further explore this later. The UNet and FSTR models show much lower RMSE in 306 water table depth across the domain, with only some regions of higher error in the western 307 portion of the domain. Overall, the FSTR model has the lowest RMSE, and highest correlations. 308 All models had relatively low correlations on the western edge and sections in the central area of 309 the domain. Interestingly, there are some differences in where the UNet and FSTR models show 310 higher errors. For instance, FSTR has a region of low correlation in the plains on the western 311 sides of Nebraska and Kansas, while the UNet seems to capture these correlations well 312 (Although the FSTR model still showed lower overall errors in this region). The model 313 performance improvement of FSTR over the UNet was primarily marked by reductions in RMSE 314 over the eastern two thirds of the domain. 315

316





318 319

320

Figure 4. Spatial metrics of performance of each model architecture's ability to emulate water table depth over the entire testing period of water year 2006. The left column (a-c) shows the root mean square error (RMSE), while the right (d-f) shows the Pearson correlation.

Similarly, in figure 5 we examine the models' ability to emulate the surface saturation levels 321 as well. Here we find again that the ResNet is the worst performing overall, with the UNet in the 322 middle, and FSTR performing best. The areas of highest error for surface moisture tend to be the 323 324 same for all the models, unlike the results for the water table depth. Of note is the northeastern portion of the domain above and around the Great Lakes, which have been masked out, as well 325 as some regions throughout the central part of the domain and in the southeast corner. These are 326 quite different hydrologic systems, which suggests deficiencies in the input data or model 327 training procedure are the underlying cause, though we were not able to diagnose the exact 328 reasons for these patterns. 329





Figure 5. Spatial metrics of performance of each model architecture's ability to emulate surface
 soil moisture over the entire testing period of water year 2006. The left column (a-c) shows the
 root mean square error (RMSE), while the right (d-f) shows the Pearson correlation.

To better understand the temporal nature of error growth of the emulators, we also show the forecast error at increasing lead times in figure 6. Here we find that the ResNet tends to have reasonable forecast error out to about day 200, before growing rapidly, leading to the pockmark error patterns from figures 4 and 5. However, both the UNet and FSTR models have significantly lower error rates and no exponential blowup over the simulation period. The FSTR model maintains the lowest errors over the entire period, like what we saw in the spatial error analysis.



#### 340

Figure 6. Growth of the error distribution of water table depth with increasing forecast length.
 The inset plot shows the full range of the error growth for the ResNet, which is more than ten
 times the error of the other two model architectures. The central lines show the median RMSE,
 while the shaded lines show the interquartile and 10-90% range of errors.

So far, we have looked at the overall error characteristics in the forecasts across the full 345 domain, but it is worth zooming in on some particular regions to see local performance spatially 346 and temporally. First, in figure 7 we show a 256 by 256 km sub-domain in the midwestern US. 347 We omit the ResNet results here because of its instability and low accuracy. Comparing the time 348 series for this region we once again see that the FSTR model is able to much more accurately 349 350 capture the dynamics of the system, as compared to the UNet. Most notably, the FSTR model is able to capture the seasonal dynamics in a much more robust way than the UNet, which shows 351 good overall correlation, but drifts in magnitude from the ParFlow results. Similarly, the spatial 352 plots show that FSTR is much more able to represent the spatial heterogeneity across seasons, 353 particularly in regions with shallow water tables. This connects back to the timeseries picture, 354 where we see the UNet consistently predicts a deeper water table, thus making it hard to form the 355 surface river network. Neither model can perfectly capture the structure of the river network, 356 although both show features that correlate well with the river network from ParFlow. 357 358







Similarly, to figure 7 we show a zoom in to the southwestern portion of the domain in figure 8. This region is much more arid and has deeper water table depths on average. Here we see that both the UNet and FSTR models can capture the long-term dynamics of the region and maintain

spatial coherence over the simulation period as well. However, we do see that the FSTR model is
 better able to capture the fine-scale structure where shallow water table depths and rivers form in
 the basins.



Figure 8. A zoom in of a region of the domain in the southwest highlighted in red on the upper

right. Timeseries shown on the top are median values for the region, while the spatial plots

372 correspond to time slices designated by dashed vertical lines in the timeseries.

In figure 7 we saw that one of the main deficiencies in the UNet output is drift in the longer-373 term dynamics, which we found in other regions as well. To better understand what drives the 374 accumulation of errors over time we looked at how the UNet and FSTR models respond to 375 precipitation events. We accomplished this by selecting grid cells at varying levels of 376 precipitation and then comparing the response of the surface level pressure heads before and 377 after the storm events (Figure 9). Overall, we found that both model architectures show similar 378 sensitivities to precipitation events, even in the extreme events. The FSTR model does show a 379 closer match to the response of ParFlow at the 90th percentile and above, which this is likely 380 contributes to the improved performance compared to the UNet. The ability of the emulators to 381 respond accurately across a wide range of events is a promising result indicating potential ability 382 to use them in evaluating the impacts of extreme precipitation events. 383



384

Figure 9. Comparison of the response of the surface layer pressure head to different precipitation
 inputs.

# 387 4 Discussion and future work

In this study we chose to focus entirely on the subsurface and treat the snow and ET as inputs to the model. We chose this because ParFlow is substantially more computationally expensive to run than the land surface model CLM, which provides the simulation of the snow, evaporation, and vegetation processes. We plan to incorporate sub-modules that simulate these processes and can be connected to the subsurface component in another study. These processes could also be added to the FSTR model, but this would require the ability to estimate their gridded initial
conditions of snowpack and ET to be used in this framework. There are also many additional
experiments that could be performed using the framework that we have developed here to further
improve the capabilities of the emulators.

One opportunity for further exploration is to modify the training routines that we use to produce more suitable emulators for specific time or spatial scales. Our aim of enabling subseasonal to seasonal groundwater prediction was our initial choice for the daily timestep that our models operate at, but we showed they can be rolled out to annual scales with little degradation in performance with our FSTR model. It may even be possible to combine training methodologies into a multi-stage model, similar to the approach taken by the FuXi weather forecasting architecture (Chen et al., 2023).

Another aspect that we have not yet explored is in the translation from the pressure heads to 404 water table depth and soil moisture during training time. The equations we used are implemented 405 in the same way as they are in ParFlow, but are written as PyTorch layers which, in theory, could 406 be appended to our model structures and trained on directly. It is possible that this could yield 407 better predictions for these quantities but may sacrifice the fidelity of the more fundamental 408 pressure heads. In future work we may explore this along with multi-objective training to capture 409 a wider range of conditions. There are some technical challenges that remain before such 410 experiments can be performed, though we believe this will be a necessary step in order to 411 412 successfully emulate the overland flow component of ParFlow as evidenced by the overall difficulty of even the FSTR model to accurately reproduce the stream network via the surface 413 soil moisture. 414

One other potential limitation of our current emulation method is that the models are dependent on the soil and vegetation classifications that were used in the original simulations. Work is ongoing to develop strategies build ParFlow ensembles with varied parameters which would provide the basis for model training to be varied in more robust ways. We hope that this will provide an even stronger model which can be used to optimize the subsurface parameter values via simulation-based inference, in turn giving better simulation results as compared to observations.

Going beyond improving the emulation of ParFlow simulations, this work enables new 422 applications for large-scale simulations. For example, having a fast and stable emulator at the 423 seasonal timescale allows for ensemble predictions, uncertainty analysis, and coupling to land 424 425 surface and atmospheric models. If future work shows that such emulation approaches can be stable over the annual-to-decadal periods this will also enable more robust studies on the effects 426 of climate change on groundwater. Additionally, our FSTR architecture should be suitable for 427 modeling in other domains where external forcings act on the state of the system, particularly 428 other areas of hydrologic and land surface modeling. Additionally, we believe that the action-429 conditioning portion of the architecture provides a natural coupling point to other model types. 430

## 431 5 Conclusions

In this study we developed a deep-learning model architecture that is a viable approach to full-system emulation of a complex spatiotemporal groundwater model at high resolution. Our results show that off-the-shelf neural network architectures like the ResNet and UNet do not have the predictive capability needed for four-dimensional hydrologic simulations. Our FSTR architecture is more accurate at reconstructing groundwater dynamics and is stable over long rollout times. Our emulators exhibit much lower computational cost compared to the original

- simulations. This opens up many new opportunities to use emulated results for ensemble 438
- 439 forecasting and model calibration through simulation-based inference.
- The FSTR model architecture developed in this study not only shows good overall 440
- performance at simulating continental-scale pressure heads, water table depth, and soil moisture 441
- but is also fast and scalable to run. As a side-benefit the model architecture is easy to 442
- conceptually understand because the model inputs directly correspond to the input data types 443
- used in most distributed hydrologic models. Initial conditions are used to set the model up for 444
- simulation, static physical parameters are used to modulate the time evolution of the system, and 445
- meteorologic forcings act on the internal state of the system. Currently each of these subsystems 446 are processed via convolutional layers before being fed into the core AC-ST-LSTM model, but 447
- fundamentally could be any neural network component. 448
- The use of deep learning for geophysical modeling is still a rapidly developing field where 449
- most applications for land and subsurface modeling are either point-scale timeseries models or 450 static spatially distributed models. In this work we have demonstrated a modeling approach that
- 451
- is not only fully spatiotemporal, but also can be used to represent multiple processes 452
- simultaneously. We consider this work to be a step in moving towards more advanced 453
- 454 representations of the subsurface, which is currently lacking compared to capabilities for weather
- and atmospheric predictions. 455

#### 456 Acknowledgments

- This research was funded by the US National Science Foundation Convergence Accelerator 457
- Program, Grant No. CA-2040542. The authors would also like to thank Bill Hasling, Amy 458
- Defnet, Amy Johnson, and Will Lytle for their assistance in developing software to deploy our 459
- models to https://hydrogen.princeton.edu/. 460

#### **Open Research** 461

- The code to train the emulators and produce the figures associated with this manuscript can be 462
- found at: https://github.com/HydroFrame-ML/hydrogen-emulator-configurable 463
- Additionally, the raw datasets used to train the emulators can be accessed via this python 464
- package: https://github.com/hydroframe/hf\_hydrodata 465

#### **References** 466

- Astrid, P., Weiland, S., Willcox, K., & Backx, T. (2008). Missing Point Estimation in Models 467 Described by Proper Orthogonal Decomposition. IEEE Transactions on Automatic 468 Control, 53(10), 2237–2251. https://doi.org/10.1109/TAC.2008.2006102 469
- Bennett, A., & Nijssen, B. (2021). Deep Learned Process Parameterizations Provide Better 470 Representations of Turbulent Heat Fluxes in Hydrologic Models. Water Resources 471
- Research, 57(5), e2020WR029328. https://doi.org/10.1029/2020WR029328 472
- Blöschl, G., Bierkens, M. F. P., Chambel, A., Cudennec, C., Destouni, G., Fiori, A., et al. (2019). 473 Twenty-three unsolved problems in hydrology (UPH) – a community perspective. 474
- Hydrological Sciences Journal, 64(10), 1141–1158. 475
- https://doi.org/10.1080/02626667.2019.1620507 476
- Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and 477
- Stabilizing Machine-Learning Parametrizations of Convection. Journal of the 478 479
  - Atmospheric Sciences, 77(12), 4357–4375. https://doi.org/10.1175/JAS-D-20-0082.1

480	Chen, L., Zhong, X., Zhang, F., Cheng, Y., Xu, Y., Qi, Y., & Li, H. (2023, June 27). FuXi: A
481	cascade machine learning forecasting system for 15-day global weather forecast. arXiv.
482	Retrieved from http://arxiv.org/abs/2306.12873
483	Cheng, Y., Musselman, K. N., Swenson, S., Lawrence, D., Hamman, J., Dagon, K., et al. (2023).
484	Moving Land Models Toward More Actionable Science: A Novel Application of the
485	Community Terrestrial Systems Model Across Alaska and the Yukon River Basin. Water
486	Resources Research, 59(1), e2022WR032204. https://doi.org/10.1029/2022WR032204
487	Condon, L. E., & Maxwell, R. M. (2019). Simulating the sensitivity of evapotranspiration and
488	streamflow to large-scale groundwater depletion. Science Advances, 5(6), eaav4574.
489	https://doi.org/10.1126/sciadv.aav4574
490	De la Fuente, L. A., Gupta, H. V., & Condon, L. E. (2023). Toward a Multi-Representational
491	Approach to Prediction and Understanding, in Support of Discovery in Hydrology. Water
492	Resources Research, 59(1), e2021WR031548. https://doi.org/10.1029/2021WR031548
493	Doury, A., Somot, S., Gadat, S., Ribes, A., & Corre, L. (2023). Regional climate model emulator
494	based on deep learning: concept and first evaluation of a novel hybrid downscaling
495	approach. Climate Dynamics, 60(5), 1751–1779. https://doi.org/10.1007/s00382-022-
496	06343-9
497	Fan, Y. (2015). Groundwater in the Earth's critical zone: Relevance to large-scale patterns and
498	processes. Water Resources Research, 51(5), 3052–3069.
499	https://doi.org/10.1002/2015WR017037
500	Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., & Hochreiter, S. (2021). Rainfall-runoff
501	prediction at multiple timescales with a single Long Short-Term Memory network.
502	Hydrology and Earth System Sciences, 25(4), 2045–2062. https://doi.org/10.5194/hess-
503	25-2045-2021
504	van Genuchten, M. Th. (1980). A Closed-form Equation for Predicting the Hydraulic
505	Conductivity of Unsaturated Soils. Soil Science Society of America Journal, 44(5), 892-
506	898. https://doi.org/10.2136/sssaj1980.03615995004400050002x
507	Haber, E., & Ruthotto, L. (2018). Stable Architectures for Deep Neural Networks. Inverse
508	Problems, 34(1), 014004. https://doi.org/10.1088/1361-6420/aa9a90
509	He, K., Zhang, X., Ren, S., & Sun, J. (2015, December 10). Deep Residual Learning for Image
510	Recognition. arXiv. Retrieved from http://arxiv.org/abs/1512.03385
511	Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. Neural Computation, 9(8),
512	1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735
513	Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects.
514	Science, 349(6245), 255–260. https://doi.org/10.1126/science.aaa8415
515	Kasim, M. F., Watson-Parris, D., Deaconu, L., Oliver, S., Hatfield, P., Froula, D. H., et al.
516	(2021). Building high accuracy emulators for scientific simulations with deep neural
517	architecture search. Machine Learning: Science and Technology, 3(1), 015013.
518	https://doi.org/10.1088/2632-2153/ac3ffa
519	Keisler, R. (2022, February 15). Forecasting Global Weather with Graph Neural Networks.
520	arXiv. Retrieved from http://arxiv.org/abs/2202.07575
521	Khan, S., Naseer, M., Hayat, M., Zamir, S. W., Khan, F. S., & Shah, M. (2022). Transformers in
522	Vision: A Survey. ACM Computing Surveys, 54(10s), 200:1-200:41.
523	https://doi.org/10.1145/3505244

524	Kochkov, D., Smith, J. A., Alieva, A., Wang, Q., Brenner, M. P., & Hoyer, S. (2021). Machine
525	learning accelerated computational fluid dynamics. Proceedings of the National Academy
526	of Sciences, 118(21), e2101784118. https://doi.org/10.1073/pnas.2101784118
527	Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall-Runoff
528	modelling using Long-Short-Term-Memory (LSTM) networks. Hydrology and Earth
529	System Sciences Discussions, 1–26. https://doi.org/10.5194/hess-2018-247
530	Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Pritzel, A., et al.
531	(2022, December 24). GraphCast: Learning skillful medium-range global weather
532	forecasting. arXiv. Retrieved from http://arxiv.org/abs/2212.12794
533	Leonarduzzi, E., Tran, H., Bansal, V., Hull, R., De La Fuente, L., Bearup, L., et al. (2022).
534	Training machine learning with physics-based simulations to predict 2D soil moisture
535	fields in changing climate. Frontiers in Water, 4. Retrieved from
536	https://www.frontiersin.org/articles/10.3389/frwa.2022.927113
537	Maxwell, R. M., Condon, L. E., & Kollet, S. J. (2015). A high-resolution simulation of
538	groundwater and surface water over most of the continental US with the integrated
539	hydrologic model ParFlow v3. Geoscientific Model Development, 8(3), 923-937.
540	https://doi.org/10.5194/gmd-8-923-2015
541	Maxwell, Reed M., & Condon, L. E. (2016). Connections between groundwater flow and
542	transpiration partitioning. Science, 353(6297), 377–380.
543	https://doi.org/10.1126/science.aaf7891
544	Maxwell, Reed M., Condon, L. E., & Melchior, P. (2021). A Physics-Informed, Machine
545	Learning Emulator of a 2D Surface Water Model: What Temporal Networks and
546	Simulation-Based Inference Can Help Us Learn about Hydrologic Processes. Water,
547	13(24), 3633. https://doi.org/10.3390/w13243633
548	O'Neill, M. M. F., Tijerina, D. T., Condon, L. E., & Maxwell, R. M. (2021). Assessment of the
549	ParFlow–CLM CONUS 1.0 integrated hydrologic model: evaluation of hyper-resolution
550	water balance components across the contiguous United States. Geoscientific Model
551	Development, 14(12), 7223-7254. https://doi.org/10.5194/gmd-14-7223-2021
552	Ott, K., Katiyar, P., Hennig, P., & Tiemann, M. (2020). ResNet After All: Neural ODEs and
553	Their Numerical Solution. Presented at the International Conference on Learning
554	Representations. Retrieved from https://openreview.net/forum?id=HxzSxSxLOJZ
555	Razavi, S., Tolson, B. A., & Burn, D. H. (2012). Review of surrogate modeling in water
556	resources. Water Resources Research, 48(7). https://doi.org/10.1029/2011WR011527
557	Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat.
558	(2019). Deep learning and process understanding for data-driven Earth system science.
559	Nature, 566(7743), 195-204. https://doi.org/10.1038/s41586-019-0912-1
560	Richards, L. A. (1931). CAPILLARY CONDUCTION OF LIQUIDS THROUGH POROUS
561	MEDIUMS. <i>Physics</i> , 1(5), 318–333. https://doi.org/10.1063/1.1745010
562	Ronneberger, O., Fischer, P., & Brox, T. (2015, May 18). U-Net: Convolutional Networks for
563	Biomedical Image Segmentation. arXiv. Retrieved from http://arxiv.org/abs/1505.04597
564	Serifi, A., Günther, T., & Ban, N. (2021). Spatio-Temporal Downscaling of Climate Data Using
565	Convolutional and Error-Predicting Neural Networks. Frontiers in Climate, 3. Retrieved
566	from https://www.frontiersin.org/articles/10.3389/fclim.2021.656479
567	Shi, X., Chen, Z., Wang, H., Yeung, DY., Wong, W., & Woo, W. (2015). Convolutional LSTM
568	Network: A Machine Learning Approach for Precipitation Nowcasting.
569	arXiv:1506.04214 [Cs]. Retrieved from http://arxiv.org/abs/1506.04214

570 Tran, H., Leonarduzzi, E., De la Fuente, L., Hull, R. B., Bansal, V., Chennault, C., et al. (2021). 571 Development of a Deep Learning Emulator for a Distributed Groundwater-Surface Water Model: ParFlow-ML. Water, 13(23), 3393. https://doi.org/10.3390/w13233393 572 573 Wang, C., Duan, Q., Gong, W., Ye, A., Di, Z., & Miao, C. (2014). An evaluation of adaptive surrogate modeling based optimization with two benchmark problems. Environmental 574 Modelling & Software, 60, 167-179. https://doi.org/10.1016/j.envsoft.2014.05.026 575 Wang, Y., Long, M., Wang, J., Gao, Z., & Yu, P. S. (2017). PredRNN: Recurrent Neural 576 Networks for Predictive Learning using Spatiotemporal LSTMs. In Advances in Neural 577 Information Processing Systems (Vol. 30). Curran Associates, Inc. Retrieved from 578 https://papers.nips.cc/paper\_files/paper/2017/hash/e5f6ad6ce374177eef023bf5d0c018b6-579 580 Abstract.html Xu, G., Zhu, X., Fu, D., Dong, J., & Xiao, X. (2017). Automatic land cover classification of geo-581 tagged field photos by deep learning. Environmental Modelling & Software, 91, 127–134. 582 https://doi.org/10.1016/j.envsoft.2017.02.004 583 Yuan, X., Shi, J., & Gu, L. (2021). A review of deep learning methods for semantic 584 segmentation of remote sensing imagery. Expert Systems with Applications, 169, 114417. 585 https://doi.org/10.1016/j.eswa.2020.114417 586 587