

# Water Resources Research®



## RESEARCH ARTICLE

10.1029/2023WR035337

### Special Section:

Advancing flood characterization, modeling, and communication

### Key Points:

- A novel differentiable routing model can learn effective river routing parameterization, recovering channel roughness in synthetic runs
- With short periods of real training data, we can improve streamflow in large rivers compared to models not considering routing
- For basins >2,000 km<sup>2</sup>, our framework outperformed deep learning models that assume homogeneity, despite bias in the runoff forcings

### Correspondence to:

C. Shen,  
cshen@engr.psu.edu

### Citation:

Bindas, T., Tsai, W.-P., Liu, J., Rahmani, F., Feng, D., Bian, Y., et al. (2024). Improving river routing using a differentiable Muskingum-Cunge model and physics-informed machine learning. *Water Resources Research*, 60, e2023WR035337. <https://doi.org/10.1029/2023WR035337>

Received 18 MAY 2023

Accepted 4 DEC 2023

© 2024. The Authors.

This is an open access article under the terms of the [Creative Commons Attribution License](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

## Improving River Routing Using a Differentiable Muskingum-Cunge Model and Physics-Informed Machine Learning

Tadd Bindas<sup>1</sup> , Wen-Ping Tsai<sup>2</sup> , Jiangtao Liu<sup>1</sup> , Farshid Rahmani<sup>1</sup> , Dapeng Feng<sup>1</sup> , Yuchen Bian<sup>3</sup> , Kathryn Lawson<sup>1</sup> , and Chaopeng Shen<sup>1</sup>
<sup>1</sup>Civil and Environmental Engineering, The Pennsylvania State University, University Park, PA, USA, <sup>2</sup>Hydraulic and Ocean Engineering, National Cheng Kung University, Tainan City, ROC, <sup>3</sup>Amazon Search, Palo Alto, CA, USA

**Abstract** Recently, rainfall-runoff simulations in small headwater basins have been improved by methodological advances such as deep neural networks (NNs) and hybrid physics-NN models—particularly, a genre called differentiable modeling that intermingles NNs with physics to learn relationships between variables. However, hydrologic routing simulations, necessary for simulating floods in stem rivers downstream of large heterogeneous basins, had not yet benefited from these advances and it was unclear if the routing process could be improved via coupled NNs. We present a novel differentiable routing method ( $\delta$ MC-Juniata-hydroDL2) that mimics the classical Muskingum-Cunge routing model over a river network but embeds an NN to infer parameterizations for Manning's roughness ( $n$ ) and channel geometries from raw reach-scale attributes like catchment areas and sinuosity. The NN was trained solely on downstream hydrographs. Synthetic experiments show that while the channel geometry parameter was unidentifiable,  $n$  can be identified with moderate precision. With real-world data, the trained differentiable routing model produced more accurate long-term routing results for both the training gage and untrained inner gages for larger subbasins (>2,000 km<sup>2</sup>) than either a machine learning model assuming homogeneity, or simply using the sum of runoff from subbasins. The  $n$  parameterization trained on short periods gave high performance in other periods, despite significant errors in runoff inputs. The learned  $n$  pattern was consistent with literature expectations, demonstrating the framework's potential for knowledge discovery, but the absolute values can vary depending on training periods. The trained  $n$  parameterization can be coupled with traditional models to improve national-scale hydrologic flood simulations.

## 1. Introduction

Riverine floods pose a major risk to human safety and infrastructure (Douben, 2006; François et al., 2019; International Panel on Climate Change (IPCC), 2012; Koks & Thissen, 2016) and are linked to stream channel characteristics. Riverine floods along large stem rivers occur when the peak flow rate exceeds the stem river conveyance capacity. The timing of flood convergence and peak flood rates are influenced by the channel's geometries and flow resistance properties (Candela et al., 2005; Kalyanapu et al., 2009). In recent years, we have witnessed many deadly riverine floods, for example, in the Mississippi River, USA (Rice, 2019) and in India (France-Presse, 2022), with such disasters expected to rise significantly based on future climate projections (Dottori et al., 2018; Prein et al., 2017; Winsemius et al., 2016). The ability to better account for flood convergence and streamflow processes is urgently needed to help us better inform society of stem river flood magnitudes and timing.

In hydrologic modeling, routing describes how the stream network conveys runoff downstream while accounting for mass balances and the speed of flood wave propagation (Mays, 2010). Most routing models are based on the principle of continuity (or mass conservation) but they differ in how the momentum equation or flow velocity is calculated. For example, the widely applied Muskingum-Cunge (MC) (Cunge, 1969; Ponce, 1986) routing method is a center-in-time finite difference solution to the continuity equation, assuming a prismatic flood wave as the constitutive relationship to simplify the momentum equation. In some other cases, the momentum equation is solved in conjunction with the continuity equation (Ji et al., 2019) with a range of simplifying assumptions, for example, ignoring inertia (Shen & Phanikumar, 2010), ignoring both inertia and pressure gradient (only slope remaining) (Mizukami et al., 2016), or including additional formulations to handle effects of scale, for example,

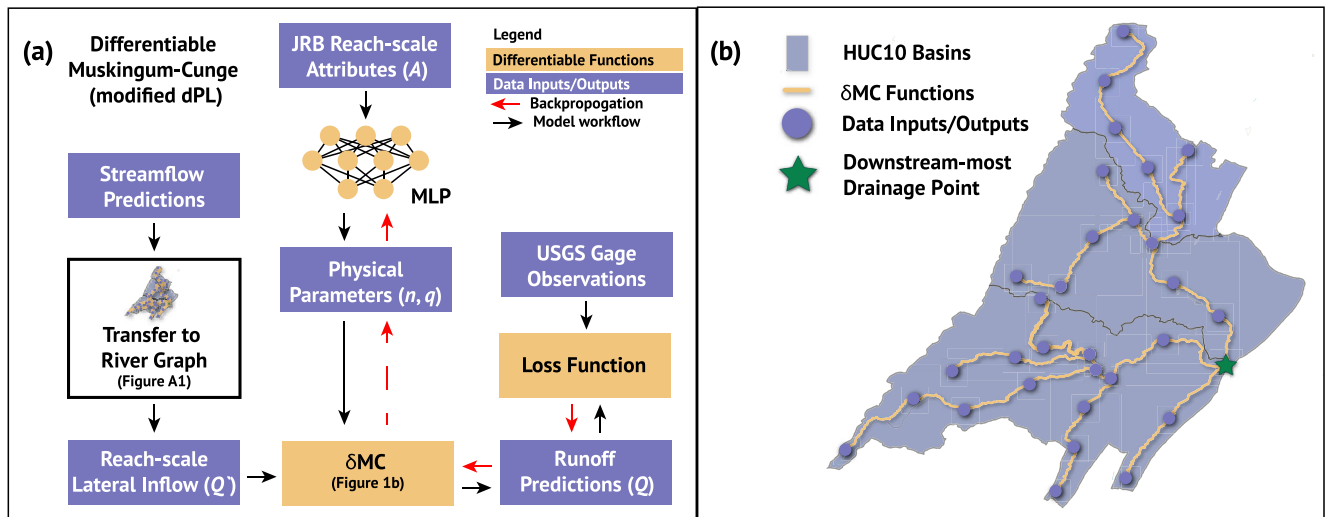
Li et al. (2013). In each case, these models have parameters that need to be determined from lookup tables or calibration, for example, roughness parameters that serve as resistance to flow.

Although routing parameters often rank among the important ones for discharge simulation (Khorashadi Zadeh et al., 2017; L. Liu et al., 2022), they have been difficult to parameterize at large scales, especially in a way to both sensibly represent basin-internal spatial heterogeneity and adapt to discharge data. Using traditional roughness values tabulated for various land covers (Arcement & Schneider, 1989) requires in situ scouting, for example, to determine if channels have pools, weeds, grass, etc., which is currently impractical for large-scale applications. Many calibration exercises (Khorashadi Zadeh et al., 2017; L. Liu et al., 2022; Mizukami et al., 2016) have used only one set of parameters for an entire basin, neglecting fine-scale spatial heterogeneity in river-reach characteristics. Some studies have employed Manning's roughness,  $n$  (a coefficient representing a channel's resistance to flow), as a linear function of river depth or other characteristics (Getirana et al., 2012; H.-Y. Li et al., 2022), but it is unclear if these relationships accurately represent the available data.

While the accuracy of basin rainfall-runoff models has improved substantially in recent years with machine learning (ML) (Adnan et al., 2021; Feng et al., 2020; Khoshkalam et al., 2023; Kratzert et al., 2019; Mangukiyi et al., 2023; Sun et al., 2022; Xiang et al., 2020), these methods have not been applied to routing modules in order to benefit the simulation of stem river floods. Neural networks (NNs) like long short-term memory (LSTM), GraphWaveNet (Sun et al., 2021), or convolutional networks (Duan et al., 2020) have demonstrated their prowess in learning hydrologic dynamics from big data. They are applicable not only to streamflow hydrology but also to variables across the entire hydrologic cycle (Shen, Chen, & Laloy, 2021; Shen & Lawson, 2021) such as soil moisture (Fang et al., 2017, 2019; J. Liu et al., 2022, 2023; O & Orth, 2021), groundwater (Wunsch et al., 2022), snow (Meyal et al., 2020), longwave radiation (Zhu et al., 2021), and water quality parameters like water temperature, dissolved oxygen, and nitrogen (He et al., 2022; Hrnjica et al., 2021; Lin et al., 2022; Rahmani, Lawson, et al., 2021; Saha et al., 2023; Zhi et al., 2021). However, these approaches are mostly suitable for relatively homogeneous headwater basins; spatial heterogeneities in forcings and basin characteristics are generally not well represented in these approaches. In our previous studies we observed that large basins often have poorer performance for LSTM models than smaller basins. The routing module is the key component that allows us to consider how runoff from heterogeneous subbasins converge and contribute to the stem river floods, and could be extended to support reactive transport modeling in the river network.

A recent development in integrating ML with physical understanding is the use of differentiable, physics-informed ML models, which can approach the performance of purely data-driven ML models but also provide interpretable fluxes and states (Shen, Appling, et al., 2023). “Differentiable” models can rapidly and accurately compute the gradients of the model outputs with respect to any input, enabling the combined training of NNs to approximate complex or unknown functions from big data while keeping physical priors (Feng et al., 2022). Such models can be simply supported by automatic differentiation (AD), which tracks each elementary operation of tensors through the use of a computational graph, then uses derivative rules to compute the gradient of each tensor operation (Baydin et al., 2018). This enables hybrid frameworks to learn and incorporate complex and potentially unknown functions from big data while retaining physical formulations. By connecting deep networks to reimplemented process-based models (or their NN surrogates), Tsai et al. (2021) developed a NN-based parameterization pipeline that infers physical parameters for process-based models. Differentiable models can also extrapolate better in space and time than purely data-driven deep networks (Feng et al., 2023). These methods are broadly applicable (including for estimation of parameters in ecosystem (Aboelyazeed et al., 2023) and stream temperature (Rahmani et al., 2023) modeling) and allow us to flexibly discover variable relationships within the model based on big data, enabling improved transparency compared to standard deep learning models.

Nevertheless, it was unclear if a differentiable model could effectively learn relationships in a highly complex river network, which convolves and integrates processes over large scales and thus renders small-scale processes unidentifiable. The river network forms a hierarchical graph, which is not unlike the graph networks for applications like social recommendations (Fan et al., 2019), but with a predefined spatial topology (due to a fixed river network) and a converging cascade. A complex river graph can have many nodes, which, when coupled with many time steps, could potentially lead to a training issue known as the vanishing gradient (Hochreiter, 1998), where the gradients with respect to the parameters are vanishingly small and the system becomes very difficult to train. Moreover, runoff data (required as an input for routing) are generally not available seamlessly for all subbasins and must be estimated by models, but models for runoff could incur substantial errors which could



**Figure 1.** (a) Flow chart of the model. Streamflow predictions are mapped to a river graph to create reach-scale lateral inflow ( $Q'$ ). Additionally, Juniata River Basin reach-scale attributes are fed into a Multilayer Perceptron network to generate physical parameters ( $n, q$ ) for  $\delta$ MC. We then predict streamflow, and evaluate it against USGS observations. (b) A visualization of how  $\delta$ MC in the river graph: the routing model is applied to each flowline within a sample river network, taking ( $n, q$ ) and the inflows of each reach as inputs, culminating in the downstream-most discharge point where loss is calculated.

potentially prevent the routing parameters from being learned. It was further unclear if downstream discharge data alone contain enough information to enable learning of reach-scale relationships. In other words, a reach-scale relationship may or may not be identifiable using downstream observations which integrate the signals from the entire catchment area.

In this work, we developed a novel differentiable modeling framework to perform routing and to learn a “parameterization scheme” (a systematic way of inferring parameters from more rudimentary information) for routing flows on the river network. The model is programmatically differentiable as it is fully implemented in PyTorch, leveraging AD. Such a physically-based routing method has never before been combined with NNs. An NN-based parameterization scheme for Manning’s  $n$  and river bathymetry shape ( $q$ ) is integrated with MC routing and is applied throughout the river network to provide improved understanding of both the model and the modeled system. We designed synthetic and real data experiments to answer the following research questions:

1. Given substantial errors with estimated runoff as inputs to the routing module, can we learn effective routing parameterization schemes that can produce reliable results for long-term simulations in large river networks?
2. Do the learned parameterization perform well for both trained and untrained downstream/upstream gages and how does performance vary as a function of basin area?
3. Do short periods of downstream discharge contain sufficient information to train a reliable parameterization scheme or to identify the parameterizations for channel roughness and hydraulic geometries?

This scheme shows promise for knowledge discovery and integration with traditional models to improve national-scale flood simulations.

## 2. Data and Methods

### 2.1. Overview

As an overview, we describe a differentiable model ( $\delta$ MC) that routes runoff through a river network (or “graph” in the ML terminology) similar to the traditional MC method. But unlike traditional MC, our differentiable model is able to track a gradient chain throughout the river routing module in order to train a connected NN that provides a reach-scale parameterization. This new routing model can be perceived as a physics-informed graph neural network (GNN) from an ML perspective. Using the gradient tracked throughout the MC routing, we trained an embedded Multilayer Perceptron (MLP) (Leshno et al., 1993) NN to generate spatially distributed river parameters for each reach (or “edge” in the GNN terminology) in the river network (Figure 1b). The training

is “end-to-end” in that there is no need for ground truth target data or pretraining for the MLP, and the whole framework is trained together based on the overall loss function. The loss function (the model's goal is to minimize the output of this equation) was calculated at the furthest downstream node of the graph. To disentangle rainfall-runoff (required information for routing) from the routing processes, lateral inflow of combined overland and groundwater flow was obtained from a pre-trained LSTM streamflow prediction model (reported in previous work and not retrained here). The runoff values were then disaggregated to hourly time steps via interpolation and routed throughout the river network using the proposed differentiable routing model (Figure 1a). We provide the details in the subsections that follow.

The full version name of the model is “ $\delta$ MC-Juniata-hydroDL2,” where  $\delta$  indicates that the model is differentiable, Juniata represents the training data set, and hydroDL2 indicates a particular software implementation (version 2.0 of the hydroDL package), but in this paper  $\delta$ MC is used as a short name for the model.

## 2.2. The River Graph

We constructed a river network (or graph) for the Juniata River Basin (JRB) in the northeastern United States (Figure 2), by processing the United States Geological Survey's (USGS's) National Hydrography Data set (NHDplus v2) (HorizonSystems, 2016; Moore & Dewald, 2016) which provides topology and some attributes of the river reaches such as upstream catchment area. We discretized the river network into approximately 2-km reaches, resulting in 544 junction points (or nodes) and 582 river reaches (or edges). These reaches are where the physical parameters like Manning's roughness and channel shape coefficients are defined. To reduce computational demand, we selected a subset of NHDplus v2 river reaches based on a stream density threshold (total stream length/watershed area), choosing rivers with the longest length until a stream density of 0.2 km/km<sup>2</sup> was reached. We then calculated slope and sinuosity for the reaches by overlaying NHDplus v2 with the 10-m resolution National Elevation Data set (USGS ScienceBase-Catalog, 2022). Prior work describes the bulk of the extraction procedure that prepares input data for a physically based surface-subsurface processes model (Ji et al., 2019; Shen & Phanikumar, 2010; Shen et al., 2013, 2014, 2016).

The hydrograph at the furthest downstream JRB gage, USGS gage 01563500 (node 4809 in our graph) on the Juniata River at Mapleton Depot, PA, was chosen as the training target (Figure 2a). This gage has a catchment area of 5,212 km<sup>2</sup> contributed from the 582 simulated reaches upstream. Seven USGS gages are located upstream of this node which enables further validation of the simulations.

## 2.3. Implementing River Routing With Muskingum-Cunge

### 2.3.1. Muskingum-Cunge

The MC method is a widely used flood routing technique that combines the Muskingum storage routing concept with the continuity and momentum equation for a river reach (Cunge, 1969), solved using a finite difference scheme for each reach, at time steps  $t$  and  $t + 1$ :

$$Q_{t+1} = c_1 I_{t+1} + c_2 I_t + c_3 Q_t + c_4 Q' \quad (1)$$

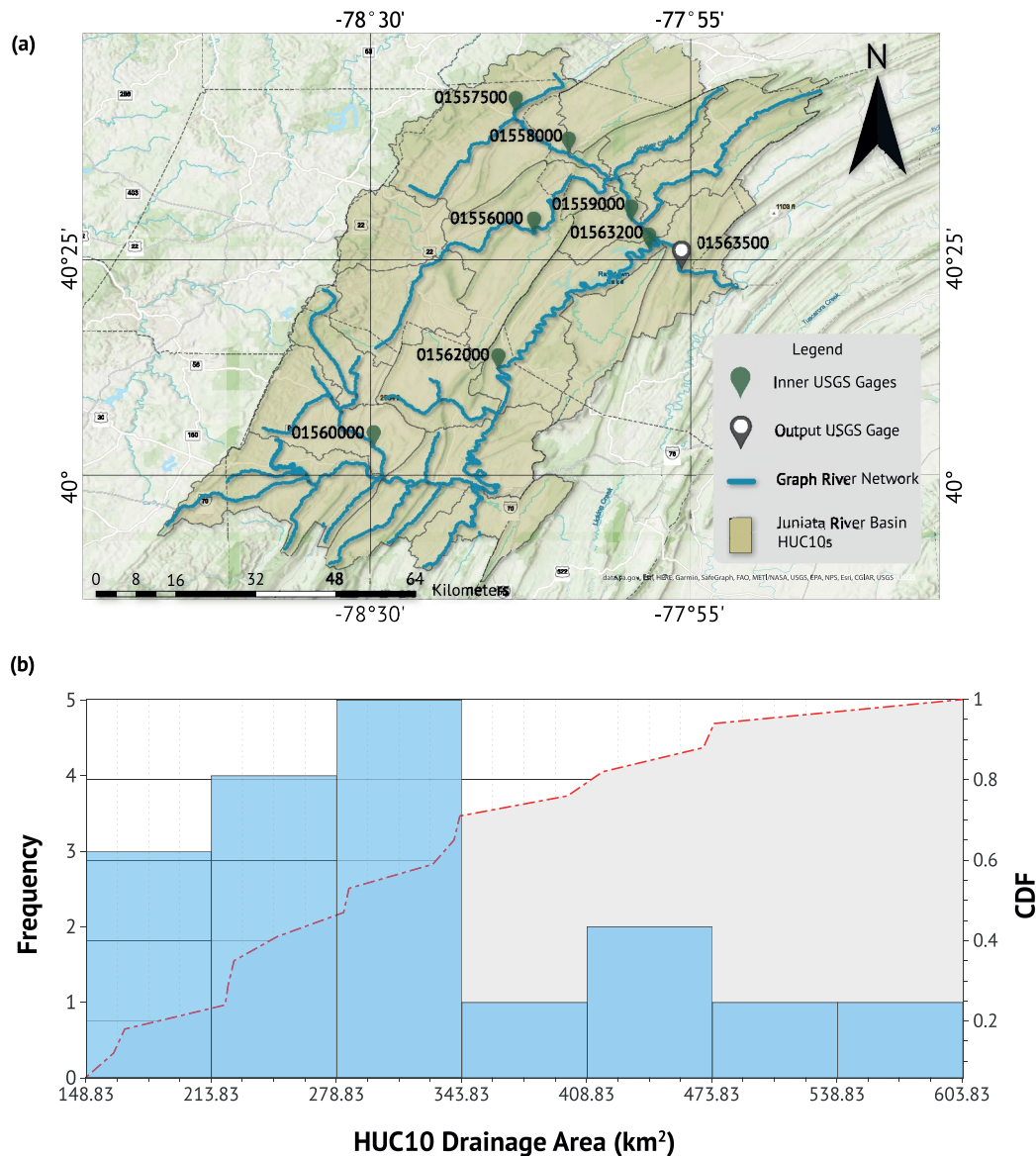
$$c_1 = \frac{\Delta t - 2KX}{2K(1 - X) + \Delta t} \quad (2)$$

$$c_2 = \frac{\Delta t + 2KX}{2K(1 - X) + \Delta t} \quad (3)$$

$$c_3 = \frac{2K(1 - X) - \Delta t}{2K(1 - X) + \Delta t} \quad (4)$$

$$c_4 = \frac{2\Delta t}{2K(1 - X) + \Delta t} \quad (5)$$

where  $I_t$  and  $Q_t$  are the inflow and outflow of the reach at time step  $t$ , respectively, and  $I_{t+1}$  and  $Q_{t+1}$  are the inflow and outflow at the next time step,  $t + 1$ .  $K$  represents travel time based on reach length and wave celerity ( $K = \Delta x/c$ ),  $X$  is a dimensionless inflow/outflow weighing parameter, and  $Q'$  represents lateral inflow of the



**Figure 2.** (a) A map of the Juniata River Basin's (JRB's) river network and HUC10 watersheds. Each eight-digit number corresponds to a USGS gage. (b) A histogram showing the distribution of HUC10 watersheds in the JRB. The x-axis shows the distribution of the HUC10 watershed area in square kilometers. The left y-axis shows the number of HUC10s that fall within the area ranges (corresponding with the blue bars), and the right y-axis shows a cumulative density function distribution of the areas, corresponding with the red dashed line.

incremental catchment area of the reach, and can also include tributary inflows (Natural Environment Research Council, 1975).

### 2.3.2. MC Parameter Values and Variable Channel Dimensions

We adopted the simple linear form of the Muskingum equation:  $X$  is constant and  $K = \Delta x/c$  where  $\Delta x$  is length of the reach and  $c$  is the celerity (m/s) of the current time step. More complex nonlinear forms of the MC equation could be tested in the future (Mays, 2019). To implement MC, we chose an hourly time step ( $\Delta t$ ) and a weighing coefficient ( $X$ ) of 0.3, which was based on regional expectations, for Equations 2–5. Since  $c$  and stream top width  $w$  vary over time, they need to be updated in each time step with respect to discharge  $Q$ , which was done here with the help of a constitutive relationship used to close the equations. For this, because at-a-site hydraulic geometries



(Gleason, 2015; Leopold & Maddock, 1953; Orlandini & Rosso, 1998) lead to a power-law relation between top width ( $w$  [m]) and depth ( $d$  [m]), we can assume such a relationship:

$$w = pd^q \quad (6)$$

where  $p$  [m] and  $q$  [–] are linear and exponential parameters, respectively, that are potentially spatially heterogeneous and represent the shape of the channel's cross-sectional area. For a rectangular channel,  $q = 0$ , and for a triangular channel,  $q = 1$ . The cross-sectional area  $A_{CS}$  is the integral of  $w$ , width, with respect to  $d$ , depth (Equation 7). To simplify the task (and because it is not sensitive based on our observations), we assumed  $p = 21$  based on preliminary data fitting to USGS hydraulic geometries from field surveys of gages in the JRB. Note that even though we make this assumption here for model completeness, we do not posit that  $q$  is invertible from available data because it may not be that significant for the downstream discharge. Moving forward with these assumptions, we can write these relationships as Equation 7:

$$A_{CS} = \int_0^d w \partial d = \int_0^d pd^q \partial d = \frac{pd^{q+1}}{q+1} \quad (7)$$

Manning's equation is:

$$v = \frac{k}{n} R_h^{2/3} S_0^{1/2} \quad (8a)$$

where  $k$  is a conversion factor equal to 1 for metric units,  $S_0$  represents the reach slope, and  $R_h$  is the hydraulic radius, equivalent to cross-sectional area divided by perimeter (Harmesen, 1955). For wide river channels, the effect of the channel sides is negligible, so depth can be used instead of hydraulic radius.

We then combine Equation 7 with Manning's  $n$  Equation and the discharge formula

$$Q_t = vA_{CS} \quad (8b)$$

where  $Q_t$  represents the discharge exiting the reach at time  $t$ . Using  $d \approx R$  to simplify the calculation, we arrive at Equation 8c:

$$Q = vA_{CS} = \frac{1}{n} R^{2/3} S_0^{1/2} \frac{pd^{q+1}}{q+1} \approx \frac{pd^{q+\frac{5}{3}} S_0^{\frac{1}{2}}}{n(q+1)} \quad (8c)$$

Reorganizing, we derive a function that estimates  $d$  from  $Q$  (Equation 8b).

$$d = \left[ \frac{Q_t n (q+1)}{p S_0^{\frac{1}{2}}} \right]^{\frac{3}{5+3q}} \quad (8d)$$

As discussed above, since channel depth can be substituted for hydraulic radius in the case of wide channels, we can use  $d$ ,  $n$ ,  $p$ ,  $S_0$ , and  $q$ , to estimate  $v$ , and then wave celerity,  $c = 5/3v$ , assuming a simple rectangular channel for  $c$ . This simplifying assumption, often made in previous large-scale routing models, is used to maintain the linearity of the equation for an easier and more stable numerical solution, as preliminary studies show that deriving  $c$  using the full equation leads to instability.

### 2.3.3. Differentiable Modeling

By implementing MC on a differentiable coding platform (PyTorch, Tensorflow, Julia, JAX, etc.), we can train a coupled NN in an “online” or “end-to-end” way to produce physical reach-scale river parameters to directly input into the routing model, much like our earlier work in differentiable parameter learning (Tsai et al., 2021) (Figure 1a). “Online” or “end-to-end” here means the routing component and the NN component are fully trained together in one stage, and no ground truth is needed for the outputs of the NN to supervise the training (although such information can be included as additional constraints when available). Rather, the training signal is back-propagated through the physical model based on what it simulates. We include an NN into the MC routing

framework to optimize equation parameters based on big data while maintaining physical consistency and mass balances. In this case, an MLP is incorporated, featuring two hidden layers and a sigmoid activation function in the output layer. The MLP accepts a normalized array of attributes ( $A$ ) for each reach (Table A2). Based on initial results, we saw no need to add further complexity (additional hidden layers), as more additions to the MLP increase the chance of overparameterization. More complexity can be added in future work when more data at larger scales are used in training and thus provide more information. The network then outputs the Manning's roughness coefficient ( $n$ ) and channel bathymetry shape coefficient ( $q$ ):

$$n, q = \text{NN}(A_{\text{JRB}}) \quad (9)$$

where  $n$  represents a channel's resistance to flow and  $q$  represents the shape of the channel's cross-sectional area. These parameters are inferred for each reach using the attributes of that reach prior to routing, since we assumed  $n$  and  $q$  to be time-invariant. This produces  $r$  number of  $n$  and  $q$  values specific to each reach for all timesteps where  $r$  is the number of river reaches. The weights of the MLP are initialized using Xavier (Glorot) initialization assuming a normal distribution (Glorot & Bengio, 2010) and are updated using backpropagation and the Adam optimizer (Kingma & Ba, 2017).

Parameters were chosen for this study to examine the relationship between channel roughness and shape, and calculated runoff. Previous studies (David et al., 2013; David, Habets, et al., 2011; David, Maidment, et al., 2011) have shown  $K$  to be more sensitive than  $X$  with respect to runoff. Thus, we believe that there is an opportunity to examine  $n$  and  $q$ , as these parameters heavily influence the calculated value of  $K$ .

#### 2.4. Lateral Streamflow Inputs

Since spatially distributed runoff is needed to predict runoff in downstream basins, but there is no such data set, we employed a pretrained LSTM (Hochreiter & Schmidhuber, 1997) rainfall-runoff model. This LSTM model was similar to those developed and reported in previous streamflow and water quality studies (Feng et al., 2020; Ouyang et al., 2021; Rahmani, Lawson, et al., 2021; Rahmani, Shen, et al., 2021), and we refer the reader to these publications for a more detailed description of these models. After the initial training was done, we chose not to further update the LSTM in order to disentangle the rainfall-runoff and routing parts of the modeling process, testing the robustness of the methodology in the face of errors with simulated runoff. In addition, the test could tell us if other rainfall-runoff models could be used instead. Updating LSTM further could lead to its co-adaptation with the routing module, making the procedure complex.

To briefly summarize, the LSTM model used a combination of basin-averaged attributes, daily meteorological forcings, and volumetric streamflow observations as inputs, and output daily basin discharge. Meteorological forcings (total annual precipitation, downward long-wave radiation flux, downward short-wave radiation flux, pressure, temperature) were obtained from the NASA NLDAS-2 forcing data set (Xia et al., 2009, 2012). We selected 29 basin attributes (Table A1 in the Appendix A) similar to those chosen in previous LSTM studies (Ouyang et al., 2021). Consistent with Ouyang et al. (2021), we focused on training the LSTM on 3,213 gages selected from the USGS Geospatial Attributes of Gages for Evaluating Streamflow II (GAGES-II) data set (Falcone, 2011) with input data between 1990/01/01–1999/12/31. We developed the workflow to obtain forcing data and inputs seamlessly for any small basin in the conterminous United States (CONUS). In this case, we extracted data from HUC8 subbasins and HUC10 watersheds to gather inputs to train our LSTM model and predict discharge, respectively.

The LSTM streamflow model achieved a median daily Nash-Sutcliffe Efficiency (NSE) (Nash & Sutcliffe, 1970) of 0.7849 for the eight gaging stations in the JRB. After training during the period of 1990/01/01–1999/12/31, the model was run from 2000/01/01–2009/12/31 to predict discharge for the 17 HUC10 watersheds in the study area:

$$Q' = \text{LSTM}(x_{\text{HUC10}}, A_{\text{HUC10}}) \quad (10)$$

where  $Q'$  [ $\text{m}^3/\text{s}$ ] is the daily runoff for the HUC10 basin, and  $x_{\text{HUC10}}$  and  $A_{\text{HUC10}}$  are HUC10-averaged atmospheric forcings and static attribute variables, respectively. Lastly, we computed a mass transfer matrix, which tabulates the fraction of a subbasin draining into a river reach. Each row of the matrix is obtained by dividing the incremental catchment area of reaches inside a subbasin by the total area of that subbasin. Runoff can be distributed to river reaches via a simple matrix multiplication (Figure A1).

Due to the nature of the data used to train the LSTM, it could produce seamless (having no gaps) runoff estimates for the JRB but only on a daily, not hourly, scale. To ensure our  $Q'$  values were on the same time-step as MC (and because MC routing needs to operate on smaller time steps), we quadratically interpolated (Virtanen et al., 2020) daily data into hourly time steps, where each daily measurement occurred at 12:00 hr. For training and evaluating the routing model, we collected observed discharge data for nodes intersecting USGS GAGES-II monitoring stations. Only some time periods of the furthest downstream gage station were used for training, and data from other stations were only used for evaluation. The observed discharge data were similarly disaggregated to hourly data.

## 2.5. Inverse-Routing and Hyperparameters

There are time zone differences between the forcing data (recorded using UTC) and USGS streamflow (recorded in UTC-5). To address this, we first shifted the LSTM-produced runoff outputs by 5 hr. Because LSTM was trained to predict runoff at the outlet of a basin, with catchment area being an impactful input to the model, it already implicitly considers the time of concentration to the outlet. However, as our modeled river network extends into the subbasins and contains smaller rivers, the routing module explicitly simulates the within-basin concentration process. Ideally, we can use an inverse-routing approach to revert LSTM-predicted runoff to the time before it enters the river network. However, as inverse-routing methods (Pan & Wood, 2013) can be quite involved and were not the focus of the study, we opted for a simple approach that shifted the runoff back in time by  $\tau$  hours.  $\tau$  is considered a hyperparameter. To avoid overfitting, we used the same  $\tau$  value for all the HUC10 subbasins and experiments and determined this value by manually tuning based on the training period. We found  $\tau = 9$  (hours) to be representative of the JRB system. More complicated procedures could be employed in the future, but this straightforward approach proved to be effective in our case.

Hyperparameters and training period sizes for our differentiable routing model were chosen through trial and error based on the training period. These trials led us to choose a hidden size of six for our MLP, and a training size of 8 weeks. Parameters converged after 50 epochs for synthetic and real data experiments. Our loss function is a combination of Mean Squared Error, a range-bound penalty, and a monotonic penalty. The range-bounded loss was implemented to ensure the MLP did not create parameters that were outside of their defined literature bounds, while a small monotonic penalty was added to low drainage area (DA) (0–500 km<sup>2</sup>) reaches to reduce the uncertainty. Since our differentiable model at  $t = 0$  assumes no inflow to the river network and relies exclusively on  $Q'$  for flow inputs, a period of 72 hr was employed to “warm up” the model states in the river network, and the loss function and NSE were not calculated within this period.

## 2.6. Experiments

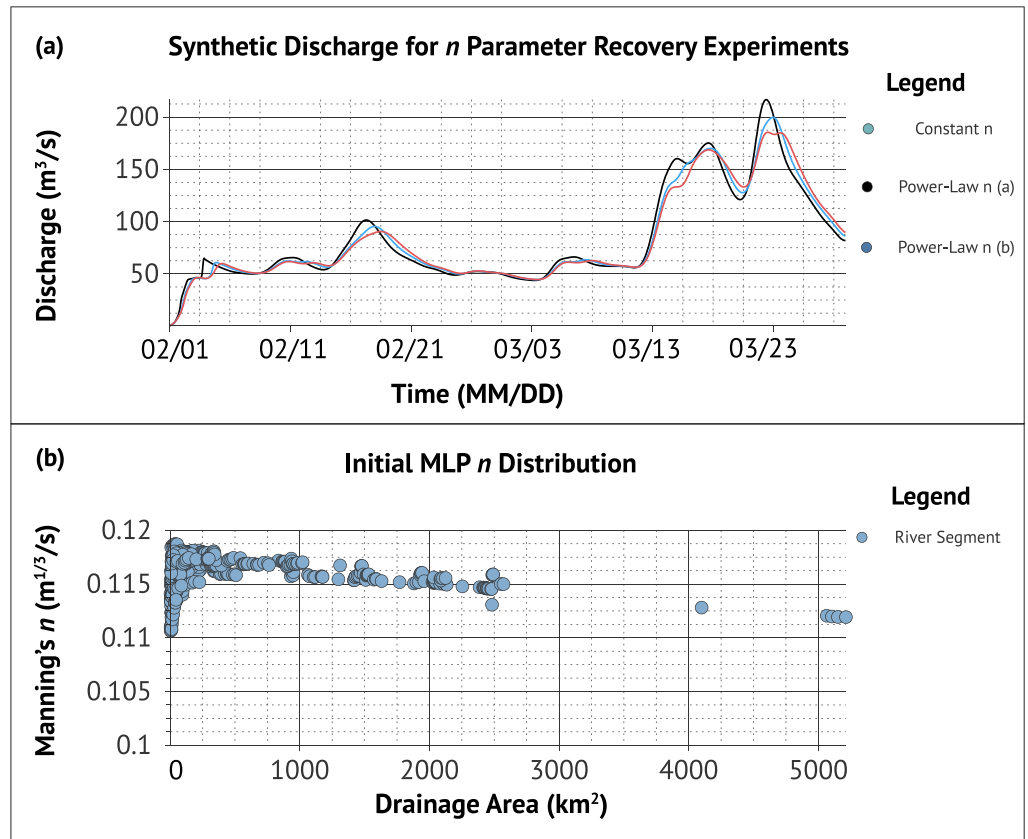
### 2.6.1. Synthetic Parameter Recovery

We first ran multiple synthetic parameter recovery experiments to check if the data set and the framework could indeed recover assumed relationships with small training periods of 8 weeks. Our first experiment tested if we could correctly recover a single, spatially constant set of assumed values for both  $n$  and  $q$  for the whole river network, resulting in only two degrees of freedom. We assumed ranges from 0.01–0.3 and 0–1 for the synthetic values of  $n$  and  $q$ , respectively, and started all of the experiment runs at the middle of the parameter range. We collectively ran 15 experiments, each with a different combination of five “synthetic truth” values of  $n$  and three “synthetic truth” values of  $q$ .

In our second experiment, we assumed constant  $n$  throughout the reaches but set the trained model as  $n, q = \text{NN}(A)$  (Equation 9) so that the  $n$  and  $q$  values could be different from reach to reach. In this case, ideally, the NN would learn to output a constant value regardless of the inputs. We set  $n = 0.08$ ,  $q = 0.5$ , and added random noise (ranged between  $-0.005$  and  $0.005$ ) to the constant  $n$  value to add heterogeneity to the experiment.

Our third synthetic experiment examined if we could retrieve simple assumed relationships within realistic literature bounds (power-law) [Equation 11] between  $n$ ,  $q$ , and DA, given that the MLP had far more inputs than just DA. The denominator's exponential coefficient  $b$  determines the bottom bound of the power-law curve, and a value of 0.131 is used for Power Law (a) while 0.05 is used for Power Law (b). The trained model still utilizes Equation 9, as we assumed we did not know the functional relationship a priori.





**Figure 3.** (a) Hydrographs showing the difference in discharge timing when using different “synthetic truth” values for  $n$ . (b) The  $n$  distribution obtained using initialized, but untrained, Multilayer Perceptron weights.

$$\begin{aligned} n &= \frac{0.0915}{(\text{DA})^b} \\ q &= \frac{2.1}{(\text{DA})^{0.357}} \end{aligned} \quad (11)$$

Each synthetic parameter recovery's unique  $n$  values are highlighted when the model output hydrographs of all three synthetic experiments are overlayed on top of one another (Figure 3a). The differences in the chosen  $n$ -values correlate with timing differences and peak discharge changes. After initialization, the MLP outputs  $n$  as shown in Figure 3b.

### 2.6.2. Observational Data Experiments

We trained  $\delta\text{MC}$  (updating the weights in the NN as in Equation 9) against observed USGS data. We utilized 8-week training periods from different years and checked whether the resulting parameters led to satisfactory routing in other years at both the training gage and untrained, inner (upstream) gages. Training periods were selected based on times when the LSTM had decent accuracy and when there was more variance into the MLP, allowing it to perform better when testing over a longer time period. Periods of such “high flashiness” in the JRB occurred during both 02/01–03/29 and 11/01–12/27, while the years 2001, 2005, 2007, and 2008 had decent LSTM accuracy, giving us eight time periods on which to train NN models. We then trained the differentiable routing models on all eight selected time periods separately to determine the sensitivity of the model performance to the selected training time period.

When interpreting model performance at inner gages, we compared results with the LSTM that modeled the whole JRB as a series of HUC10 basins, with a simple summation of the  $\tau$ -shifted LSTM runoff inputs used in  $c_4$  calculation ( $Q'$ ) (Equation 5). We also explored whether using a combination of inner gages, along with the

**Table 1**  
*Results From the Constant Synthetic  $n$  and  $q$  Parameter Recovery Experiments*

Run	$n$			$q$		
	Initial guess	Synthetic truth	Recovered parameter	Initial guess	Synthetic truth	Recovered parameter
1	0.145	0.03	0.0296	0.5	0.0	0.2890
2	0.145	0.04	0.0403	0.5	0.0	0.3072
3	0.145	0.05	0.0509	0.5	0.0	0.3305
4	0.145	0.06	0.0618	0.5	0.0	0.3492
5	0.145	0.07	0.0723	0.5	0.0	0.3737
6	0.145	0.03	0.0301	0.5	0.5	0.3491
7	0.145	0.04	0.0419	0.5	0.5	0.4016
8	0.145	0.05	0.0536	0.5	0.5	0.4517
9	0.145	0.06	0.0662	0.5	0.5	0.4809
10	0.145	0.07	0.0785	0.5	0.5	0.5211
11	0.145	0.03	0.0312	0.5	1.0	0.1930
12	0.145	0.04	0.0408	0.5	1.0	0.1750
13	0.145	0.05	0.0502	0.5	1.0	0.1649
14	0.145	0.06	0.0588	0.5	1.0	0.1801
15	0.145	0.07	0.0671	0.5	1.0	0.1983

furthest downstream gage, inside of the loss function would improve model performance on all gages throughout the study area. The gages used were USGS 01560000 (edge 1053) and 01562000 (edge 2662). Internal gages were selected based on NSE metrics when using only the furthest-downstream gage in the loss calculation; we chose basins with middle-level NSE values so as to not overfit the model if using highly predictive internal gages.

### 3. Results and Discussion

In the following, we first discuss our synthetic experiments (Section 3.1) which explore our routing framework's potential to retrieve assumed parameters from our differentiable GNN. Next, we show the results of confronting our model with LSTM-simulated runoff as observed streamflow at the furthest downstream gage, expanding the training period to other time ranges, then applying our models to different years for observation (Section 3.2). Furthermore, we discuss the stability of our trained models over several years of testing (Section 3.3). Lastly, we analyze the  $n$  parameters recovered for the trained models and discuss their implications (Section 3.4).

#### 3.1. Synthetic Experiments

Our first synthetic experiment (with constant parameters and only  $2^\circ$  of freedom for the search) recovered the assumed  $n$  values with moderate accuracy, but not the channel geometry parameter  $q$  (Table 1). Recovered  $n$  values were within a small range of the assumed ones, with minor fluctuations, while recovered  $q$  values mostly stayed similar to the initial guesses, showing

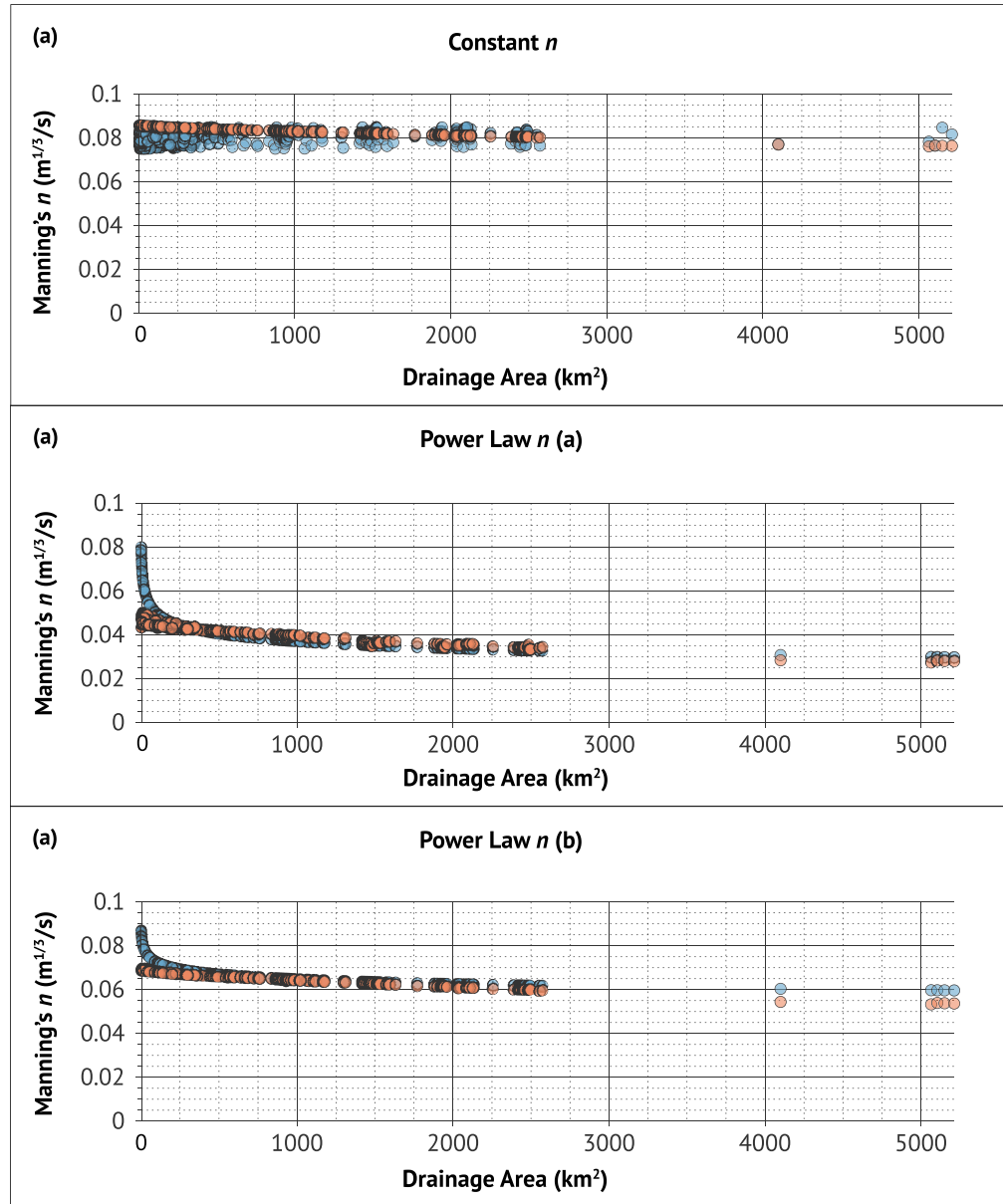
some slight changes after a number of iterations. The result was consistent across 15 runs, each starting at the same points for  $n$  and  $q$  but with a different combination of synthetic true values. The training led  $n$  to the assumed values rapidly, typically within 20 epochs (Figure A2 in Appendix A). The non-identifiability of  $q$  was likely because  $q$  has only a small influence on the storage capacity of the stream. Indeed, the time-averaged gradient of the loss function with respect to  $n$  is much larger than that with respect to  $q$  (Figure A3 in Appendix A). While it is a pity that hydraulic geometry parameters cannot be estimated, the results also implied that they would not influence the routing results noticeably. Thus, in our next efforts, we focused on  $n$ .

The second and third set of synthetic experiments together showed that either constant values or simple functions could be roughly recovered for most of the reaches, while there may have been additional uncertainty for the furthest upstream and downstream reaches (Figure 4). The constant value recovery shows retrieved values that are at the same level as the assumed range (Figure 4a). For both power-law experiments, while the two experiments differed in their mean values, in both cases the estimated  $n$  outputs from the MLP overlapped to a great extent with the value to be retrieved. The  $n$  values and the declining trends were somewhat underestimated for the headwater reaches (small-DA). In the middle ranges of DA, the curve followed the assumed one almost exactly. Toward the higher range of DA, the recovered values were lower than the assumed relationship, but the deviation was not huge because the power-law formulation became flat in this range. Based on the closeness of hydrographs in Figure 3a, we do not anticipate that further optimization can bring significant improvement to the estimations. Similar to the two-constant-parameter retrieval experiment, the  $q$  parameter was not recoverable (Figure A4 in the Appendix A).

Based on these simple experiments, it seems training on the river graphs has some promise but also some limitations. It is promising because it is likely that  $n$  is related to DA which is, to some extent, recoverable. It is simultaneously challenging because, as a large number of reaches contribute to one gage, it is an underdetermined system. This method was not able to fully reproduce the drastic change in the low-DA range presumably because this sharp slope was inconsistent with the rest of the curve, and NNs generally do not output extreme values. It also ran into difficulty toward the high-DA range because there were simply far fewer reaches with large DA so their roles in routing were relatively minor, making the curve unconstrained in this range. This experiment informed us we should not expect values of reach-scale  $n$  to be highly reliable, but the overall trend and mean

### Synthetic $n$ Values Compared to Recovered $n$ Values

● Synthetic  $n$  ● MLP recovered  $n$

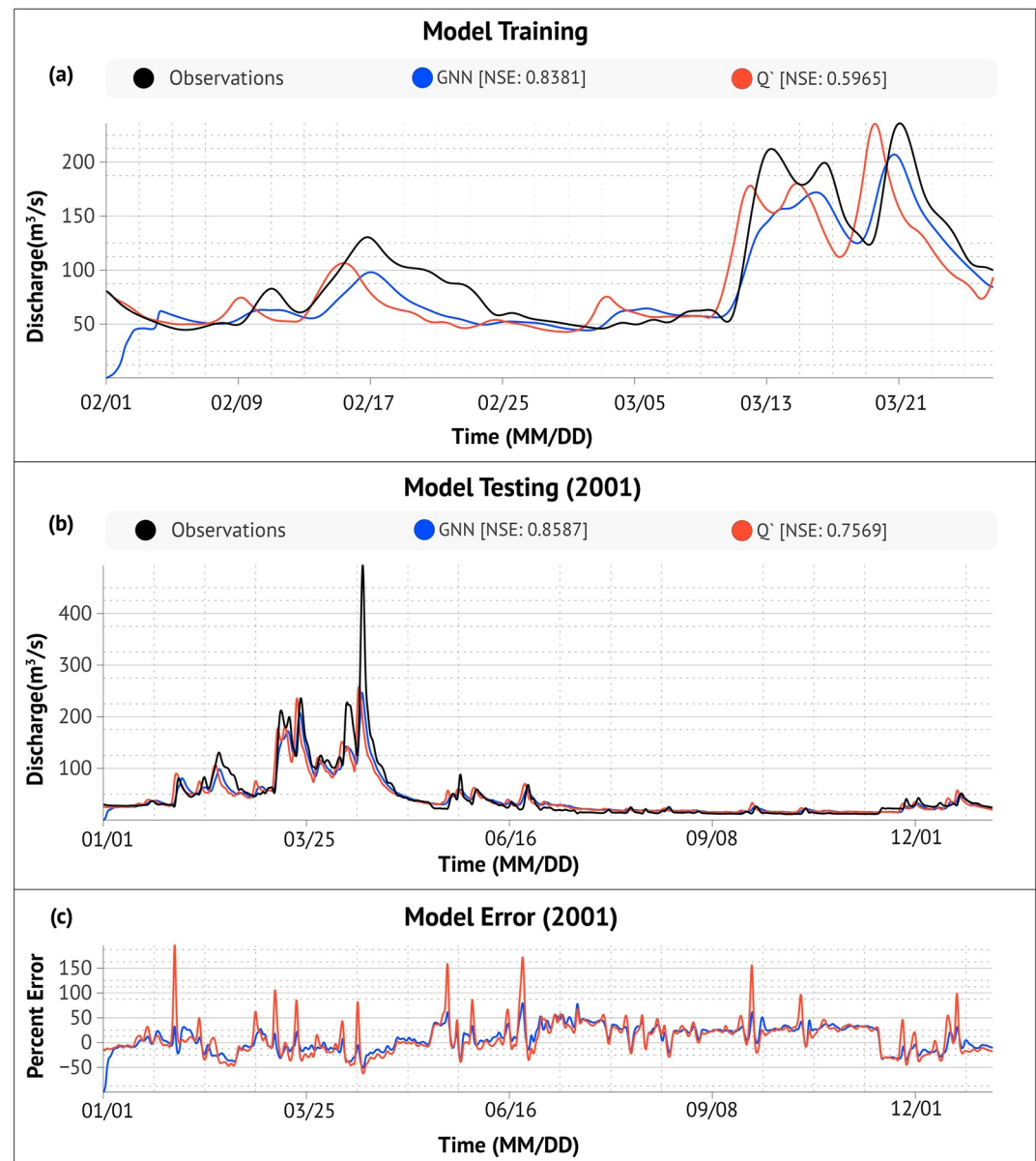


**Figure 4.** Synthetic modeled values of  $n$  with respect to the reach's total drainage area ( $\text{km}^2$ ). The neural network can recover the overall pattern but is not accurate near sharp changes or for reaches with large drainage areas. Each dot in the scatter plots represents a 2-km river reach in the river network.

values may have merit, especially when we also have other constraints. These findings formed the basis for the next stage of the work where we trained  $n = \text{NN}(A)$  for real-world data. We thus expected to extract the overall patterns of  $n$  distribution but the recovered  $q$  would not be meaningful.

### 3.2. Training on 8 Weeks of Real Data

The real-world data experiment showed satisfactory streamflow routing in the training period (Figure 5a), with improvements compared to approaches that did not employ the routing scheme, even though there was bias in the  $Q'$  term (Figure 5a). The hydrograph generated by the differentiable routing model was, as expected, smoothed and delayed compared to the summation of runoff generated during the training period. Unlike the



**Figure 5.** (a) Results from training the differentiable model (graph neural network) during an 8-week period (2001) against USGS observations compared with the summation of lateral inputs (denoted by  $Q'$ ). (b) Results from testing the trained model from Figure 4a over a year-long period (2001) compared with the summation of lateral inputs. (c) The percent error chart of both the  $Q'$  and differentiable routing outputs to show how river routing is more stable and more accurately times the flood peaks compared to using a summation of lateral inputs.

direct summation of the runoff, which has a timing difference from the observation, the peaks of the routed hydrograph are placed almost exactly under the observed peaks, leading to a high training NSE of 0.8381. We noticed a substantial low bias in this training period, witnessed by much lower peaks with the simulated flow compared to the observed flow. This is due to bias in the rainfall-runoff modeling component (arising potentially from low bias in precipitation) and the mass balance dictated by the MC formulation, which prevents the model from adding or removing mass to remove the bias. In traditional hydrologic model calibration, bias can be a significant concern as it can distort other parameters. In this case, we found the model performed well even with such bias, and appropriately focused on adjusting the timing of the flood waves. This is because the allowable adjustments were limited to routing parameters, which blocked the model from distorting other processes.

**Table 2**

*Internal Gage Nash-Sutcliffe Efficiency Values for the Year 2001, With the Rows Ranked by the Size of the Subbasin From Small to Large*

Edge ID	USGS gage number	Basin drainage area (km <sup>2</sup> )	Uniform LSTM	Q' runoff NSE ( $\tau = 9$ )	Differentiable routing model ( $\tau = 9$ )	Multiple gage loss for differentiable routing ( $\tau = 9$ )
1280	01557500	94.8	<b>0.8149</b>	0.6310	0.5611	0.5611
1053	01560000	440.5	<b>0.7028</b>	0.6129	0.6551	0.6556
2799	01558000	542.1	<b>0.8201</b>	0.7513	0.7734	0.7735
4780	01556000	723.5	0.6624	0.6615	0.6925	<b>0.6927</b>
2662	01562000	1943.5	0.7957	0.6869	0.7978	<b>0.7985</b>
4801	01559000	2103.0	0.7815	0.7468	0.8132	<b>0.8136</b>
2689	01563200	2482.9	0.5703	0.6497	<b>0.7822</b>	0.7812
4809	01563500	5212.8	0.8024	0.7569	<b>0.8587</b>	0.8583

*Note.* The differentiable routing model was trained on the period from 2001/02/01 to 2001/03/29 and the loss function was calculated for the final, furthest-downstream gage, but the LSTM was trained using >3,000 CONUS gages. We include the LSTM NSE to show how the use of routing compares to just using LSTM predictions. Bold font indicates the top performing model for each gage.

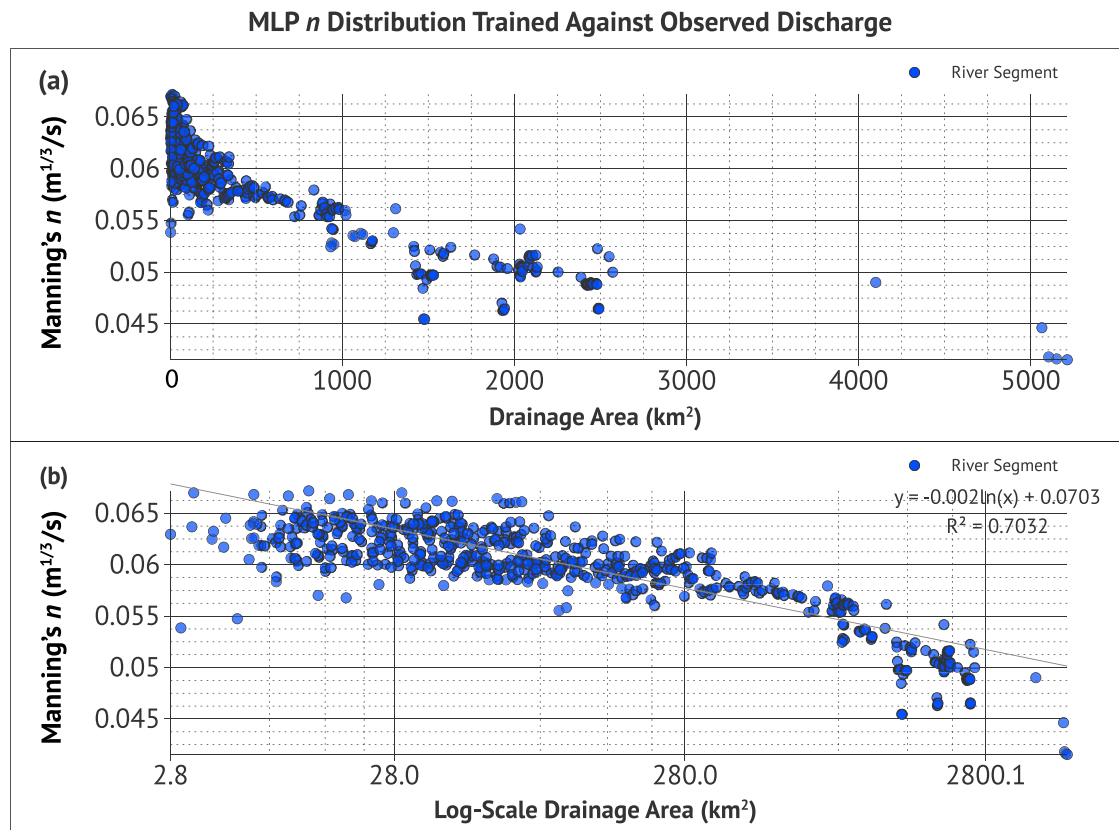
The year-long test of the differentiable model yielded high metrics compared to the alternatives (Figure 5b), suggesting a short calibration period could yield parameterizations suitable for long-term simulations. The differentiable model obtained a year-long NSE of 0.8587, which is consistent with the median NSE in the JRB. In contrast, the summation of  $Q'(\tau = 9)$  and the whole-basin LSTM were at 0.7569, and 0.801 (Table 2), respectively. This comparison shows that if we merely added the runoffs together (which already resolved spatial heterogeneity in runoff but not the flow convergence process), the error due to timing could reduce NSE at the downstream gage. While the model had success with correctly timing the peak flows, it could not compensate for LSTM's errors, resulting in significant underestimation of the peak events. By design, the routing module should be detached from the errors in the runoff module.

Interestingly, without specific instructions, the scheme recovered a power-law-like relationship between  $n$  and DA (Figure 6), similar to the one assumed in the synthetic case (Figures 4e and 4f). The  $n$  values were highest (near  $n = 0.66$ ) for smaller DA and declined gradually, approaching 0.04 at the lower end. The change rate of  $n$  as a function of DA then became more gentle as DA increased. This distribution agreed well with the general understanding (Camporese et al., 2010; Orlandini, 2002) that headwater streams running down ridges (this region is characterized by Ridge and Valley formations) have larger slopes, higher roughness, more vegetation, and thus higher  $n$ , while the high-order streams in the valley tend to have smaller slopes and smoother beds, corresponding with lower  $n$ . In most hydrologic handbooks (Mays, 2019), a smaller  $n$  is prescribed for larger rivers. Nevertheless, an advance here is the ability to train a function without a priori assumption of its form, along with the routing model, rather than having a multistep process.

### 3.3. Inner Gage Evaluation and Effects of Different Training Periods

Evaluating the model on the inner, untrained gages showed that the routing scheme became more competitive compared to benchmark levels for downstream gages (Table 2). As for the benchmarks, the uniform LSTM (the catchment area of each gage is considered a basin and basin-averaged forcing/attributes were used as inputs to the trained LSTM to simulate flow at the gage) already attempts to consider routing internally but does not consider rainfall/attribute spatial heterogeneity, while the summation of  $Q'$  (runoffs were simulated from multiple HUC10 basins and added together) considers the spatial heterogeneity but not routing in the stem river. For two of the four gages with larger than  $\sim 2,000$  km<sup>2</sup> of catchment area, the differentiable routing model performed noticeably better than the uniform LSTM models (for the other two, they were about the same). For the three mid-sized subbasins (500–2,000 km<sup>2</sup>), the comparisons were mixed. For the small subbasins, and especially gage 01557500 (94.8 km<sup>2</sup>), the uniform LSTM was noticeably better. The subbasin for 01557500 is smaller than our runoff-producing unit (HUC10s, with the smallest one  $\sim 200$  km<sup>2</sup>). This means predictions below this threshold can be error-prone. Our differentiable model was also better than the summation of  $Q'$  for seven of the eight gages





**Figure 6.** The learned relationship between  $n$  and drainage area (DA) (square kilometers) for the Juniata River basin according to the trained graph neural network. (a) The distribution of river segments by Manning's  $n$  and DA on a linear scale. (b) The same distribution, but on a logarithmic scale with a logarithmic trendline ( $R^2 = 0.7032$ ). The network was trained for the period of 2001/02/01–2001/03/29. Each dot in the scatter plot represents a 2-km river reach.

with the gap being larger for downstream gages (Table 2), suggesting the flow convergence process matters more and more as we go downstream.

When we used multiple internal gages within the NN loss function (overall loss being the sum of the losses calculated for each gage), results improved very slightly at smaller DA gages, while barely degrading at larger DA reaches. Two competing impacts are at play when we add inner gages to the training data set. On the one hand, a model calibrated at a downstream gage may perform well at the calibrated gage while producing large errors at the subbasin scale, but adding internal gages into the calibration disambiguates the contributions from tributaries and thus forces the model to be more spatially realistic. On the other hand, forcing errors tend to be larger at the subbasin scale and thus the calibration at small-DA gages may inadvertently propagate such errors and overtrain the MLP giving unrealistic  $n$  values. Overall, neither effect is strong in this case—the differences introduced by calibrating at inner gages are too small to have real-world implications, which suggests the model calibrated only at the downstream gage already had decent performance for the tributaries. The multi-gage calibration only produced a slightly more balanced model—it improves simulations at some previously weakly simulated tributaries, at a (very minor) cost to performance at the furthest downstream gage. This small tradeoff may be due to spatial errors in forcing data, as mentioned above. As the model explicitly simulates flows in all modeled reaches, the differentiable model provides a way to absorb data from as many stations as possible, if the ungauged regions are important to the users.

The above comparisons informed us of the favorable and unfavorable ranges of applicability for our workflow: the differentiable model found competitive advantages for stem rivers with catchments greater than 2,000  $\text{km}^2$ , but may run into issues for scales smaller than the smallest runoff-producing unit (HUC10, around 200  $\text{km}^2$ ). The issues for the smallest basins could be attributed to the procedure that transfers mass from HUC10 subbasins to regular grids on the river network, which should be improved in future work. As a result, the smallest headwater basins are best to be directly simulated by the uniform LSTM model. Also, smaller runoff-generating units could

**Table 3**

*The Nash-Sutcliffe Efficiency Values Correspond to Testing Differentiable Routing Models on Different Test Years*

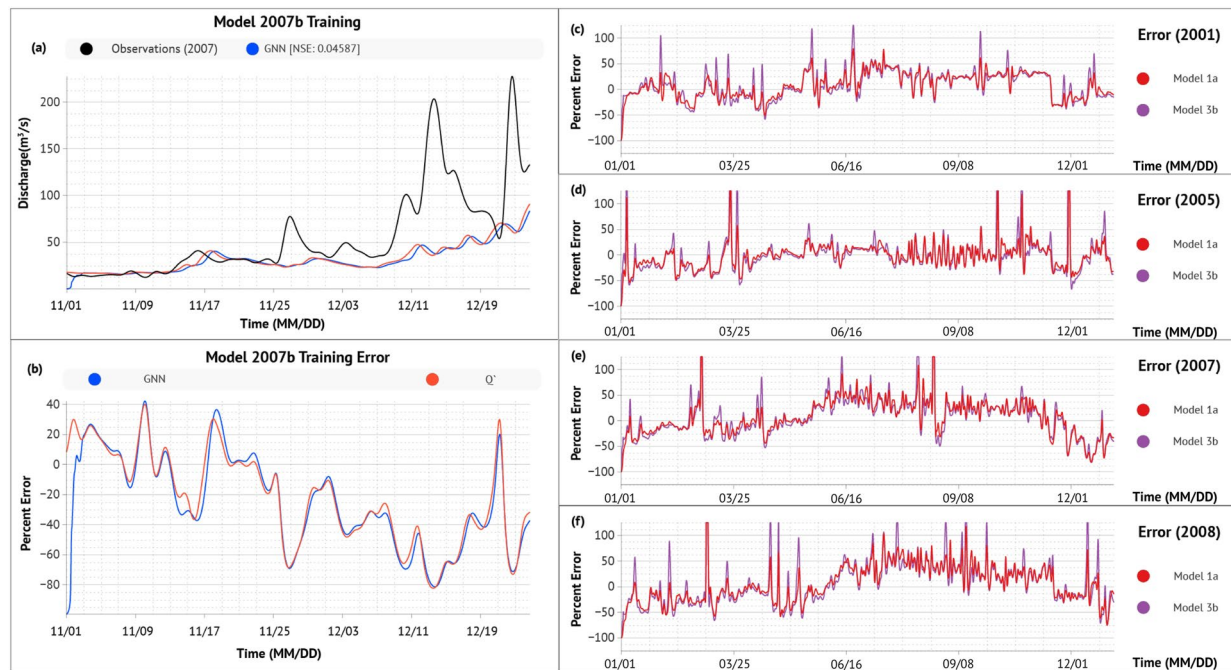
Testing period	Training period							
	2001a 02/01–3/29	2001b 11/01–12/27	2005a 02/01–3/29	2005b 11/01–12/27	2007a 02/01–3/29	2007b 11/01–12/27	2008a 02/01–3/29	2008b 11/01–12/27
2001	<b>0.859</b>	0.816	0.857	0.812	0.857	<u>0.805</u>	0.811	0.853
2005	0.813	0.817	<b>0.864</b>	0.811	0.821	<u>0.651</u>	0.810	0.761
2007	0.818	0.793	<b>0.828</b>	0.787	0.819	<u>0.738</u>	0.787	0.804
2008	0.671	0.766	0.760	0.796	0.687	<u>0.450</u>	<b>0.797</b>	0.588
Average	0.790	0.798	<b>0.827</b>	0.801	0.796	<u>0.661</u>	0.801	0.752

*Note.* Bold font indicates the highest NSE, while underlined metrics indicate the lowest NSE (noticeably worse than obtained from other periods) for the testing period.

be used in the future to mitigate this issue. As a side note, this result also indicates that LSTM-based runoff is valid down to basins as small as 94.8 km<sup>2</sup>. The advantages of the differentiable routing model over the uniform LSTM for larger basins were due to resolving the heterogeneity in rainfall and basin static attributes as well as better representing routing. The uniform LSTM can internally represent some flow lags but it appears less effective as basin size increases.

The results imply that the advantages of the differentiable routing model will increase for even larger basins where currently LSTM does not apply well, along with basins where rainfall heterogeneity makes a big difference. The JRB is situated in the northeastern part of the CONUS; many other regions may exhibit more prominent effects of heterogeneity. For example, past studies have always found it difficult to simulate large basins on the northern and central Great Plains (Feng et al., 2020; Martinez & Gupta, 2010), potentially due to spatially concentrated rainfall and runoff generation (Fang & Shen, 2017). Also, in the mountainous areas of the northwestern and southeastern CONUS, orographic precipitation could have significant spatial concentrations. We hypothesize that applying models to smaller basins and incorporating the routing scheme will allow these regions to be better modeled.

As expected, the training periods selected exerted an influence on the model, but as long as we used reasonable training periods, the results were acceptable. When the scheme was trained on 8-week periods from different years, it generated somewhat different but mostly functional parameterizations (Figure A5 in the Appendix A), unless it was trained in some unreasonable training periods where the LSTM had drastic differences from the observed outflows (Table 3). The maximum achievable NSEs for the years of 2001, 2005, 2007, and 2008 were 0.859, 0.864, 0.828, and 0.797, respectively, with all models outperforming  $Q'$  NSE values for their respective periods (Table A3 in the Appendix A). We found that if the models were trained on other periods (2001a, 2001b, 2005b, 2007a), the test NSE values were mostly decent, and at least not drastically worse. However, choosing 2007b or 2008a led to notably inferior results in both NSE values (Table 3) and in timing (Figures 7c–7f). Examining the characteristics of the different training periods, we see that the problematic training periods did not contain full flood rise and recession phases (Figure 7a) or improve timing when compared to  $Q'$  (Figure 7b). As a result, 2007b and 2008a as training periods led to either the lowest or the highest  $n$  values and also had relatively low NSE values (Figure A5 in the Appendix A). Similarly, training period 2005a gave relatively large  $n$  values which also resulted in suboptimal (although still decent) results in all the years, likely due to low  $Q'$  NSE. Hence, we conclude that periods selected for training should (a) contain full flood rise and recession phases; and (b) have high  $Q'$  NSE values. In addition, even though the routing simulation can be improved by short training periods, the spread of estimated  $n$  again shows that the identification of  $n$  via small training periods can be difficult. This is because not enough information exists in a short period to distinguish between possible solutions of  $n$ , similar to the problem of parameter equifinality or “non-uniqueness” (Beven, 2006) in rainfall-runoff modeling. We do not have sufficient events with spatially different forcings to inform the model of the regional differences. Future work needs to employ longer training periods to compromise across different periods and obtain broadly performant parameterizations. As different storm events contain varied spatial rainfall patterns and discharge magnitudes, longer training data can better constrain the model and reduce the uncertainty of the inversion. However, another possibility is that  $n$  itself can vary over time, which would be an unorthodox but not unthinkable idea.



**Figure 7.** (a) Two training periods: 2001a and 2007b trained during the time periods 2001/03/29–2001/03/29, and 2007/11/01–2007/12/27, respectively. The former contains a full rising-recession cycle (Figure 5a) while the latter does not have a complete cycle for training, thus leading to larger errors during testing. (b) Percent error of the  $Q'$  and routing outputs. (c)–(f) Percent error for the testing years 2001, 2005, 2007, and 2008, respectively. The range of the percent error charts has been clipped to a range of  $[-125\%–125\%]$  to highlight that differences in the effect trained  $n$  parameters have on flood wave timing.

### 3.4. Further Discussion

Although the estimated  $n$  values were both functional for routing streamflow and physically meaningful, the results suggest the downstream discharge only poses a moderate constraint on the  $n$  values, and short training periods may not be sufficient to identify the true  $n$  values. Hence, while our procedure can obtain an  $n$  parameterization performant for long-term simulations, we do not claim that the procedure retrieved the “true”  $n$  parameterization. Especially considering there are many input variables to the NN that covary in space, it may be difficult to disentangle causation from correlation. Due to the lack of ground truth for  $n$  in the real-data case, we leave this evaluation for future effort as we compile more measurement data. Recall that we were able to retrieve the overall pattern of  $n$  in the synthetic experiments but faced large uncertainties in some areas of the parameter space. This is attributed to the numerous degrees of freedom (a high-dimensional input space for the NN, influencing many reaches) constrained by only one downstream output with a relatively short training period. Nevertheless, this training is valuable because discharge data can be widely available, and we will be able to employ it in conjunction with other constraints, for example, scattered measurements or expert-specified relationships.

Regarding other potential recoverable parameters, we suspect (and preliminary tests show) the dimensionless MC inflow/outflow weighing parameter  $X$ , which indicates the shape of the assumed flood prism, cannot be identified for the same reason as  $q$ —the geometries of the channel do not impact flow rates in a meaningful way. While using a constant  $X$  is limiting, future work with larger data sets could reinvestigate if learning  $X$  produces any benefit. Similarly, linear channel coefficient  $p$  values were also never recoverable in single parameter tests and decreased resulting NSE values when used as a tunable parameter. We hypothesize  $p$  has a smaller effect on the width equation, and  $K$  calculation, than exponential channel coefficient  $q$ . Thus, we did not include a learnable  $p$  in the rest of the study. In addition, we hypothesize using a more complex MC formula, for example, the nonlinear form of the Cunge equation (where the celerity is defined as  $dQ/dA$ ) and a more complex treatment of the hydraulic radius, which might add to numerical challenges for large-scale simulations, would lead to different  $n$  values, as the recovered values are inherently linked to the inverse model employed.

Here we employed a static parameterization scheme for  $n$ , following the conventional approach. However, the framework allows for the use of a dynamic  $n$  (likely dependent on  $Q$ ). It is not clear if we must use a static

parameterization as done conventionally, as some previous studies have found a dynamic  $n$  to offer better results (Ye et al., 2018). In the future, it will be interesting to see if a dynamical  $n$  parameterization could significantly impact the routing results. On another note, we chose an 8-week time period as our training length as a probe to assess the required training duration and selection criteria for such periods. We trained eight different models (Section 3.3) on different time periods and showed that the choice of training period timing, and LSTM performance for the inputs, played important roles. Future effort should include longer training periods and larger basins to most robustly train the NN, and more benchmarks against alternative options such as a power-law function of catchment area or a linear function of flow depths (Getirana et al., 2012; H.-Y. Li et al., 2022), etc., can be carefully carried out. The complexity and number of the weights of the connected NN are worth further investigation as the system is evolved to learn more parameters and is confronted with more data.

When investigating the impact of multiple gages, rather than a single downstream-most gage (in model loss calculation and parameter updates), results were very similar in terms of NSE score and recovered Manning's  $n$  parameters. We believe this may be because the JRB is a relatively small river network, so internal gage observations are highly correlated in discharge volumes ( $\text{m}^3/\text{s}$ ) and fluctuations (storm event timing). Adding more gages could be useful if flows in different parts of the basin need to be accurately reported, but may be less important if only the downstream gage is of concern. We theorize increased benefits from the usage of many internal gages within the loss function as reach DA increases, due to the decoupling of gage hydrographs and storm events. Future work with large-DA studies can address this hypothesis.

We chose to parameterize channels using  $n$  and  $q$  (although the latter cannot be recovered), unlike other studies that used  $K$ , due to multiple reasons. First,  $n$  and  $q$  are here thought to be more fundamental roughness and geometric concepts than  $K$ , which is also dependent on the spatial step size of discretization. Representing the concepts at this level allows us to obtain insights such as “roughness is learnable but channel geometry is difficult to infer from downstream discharge.” That being said, using  $K$  or concepts like “floodwave attenuation and delay factors” would also produce functional models and insights. Furthermore, this work is a demonstration of the unique capability of differentiable modeling (Shen, Appling, et al., 2023), where one can freely set priors to express their understanding of the physical system and ask questions of interest to them. Here, admittedly, the values of parameters obtained via inversion, like in most other inversion tasks, would be more or less impacted by the model employed. In other words, it is very difficult to inversely obtain parameter values that are truly independent of the modeling choices. For example, using the nonlinear MC equations may generate somewhat different  $n$  values. In addition, the simplifications in Equation 8, the many minor choices such as  $\tau$ , training periods, and hyperparameters of the NNs, etc., could all potentially have minor impacts on the results. Nevertheless, we expect these values to be highly correlated to the true physical concepts and could be used to represent the spatial patterns of these parameters. Besides inversion, future work could use measurements as constraints to make the inversion more robust and physically realistic.

Our approach, akin to a classical routing scheme, is modular—the trained weights of the NN that generates  $n$  are not tied to a particular rainfall-runoff model. Our work can be coupled to traditional models in multiple ways. Firstly, the trained network can be used to generate  $n$  for traditional models. In this way, no change is required on the part of the traditional models. Secondly, the NN and the trained weights can be ported to other programming environments like Fortran. This makes it possible to use the trained parameterizations as a built-in module in continental-scale models (Greuell et al., 2015; Johnson et al., 2019; Regan et al., 2018). An alternative approach is to lump both the routing and runoff simulations into one problem and optimize them together, as demonstrated in some other studies (Jia et al., 2021). In our case, this would mean that we would train both the runoff LSTM and the routing module together. In many big-data DL case studies, lumped models tend to have higher performance compared to a workflow that separates the tasks into multiple minor tasks. However, in our case here, this leads to coadaptation concerns. Moreover, our approach is modular so it can be easily coupled to other runoff models, for example, a non-differentiable traditional model, or a differentiable one (Feng et al., 2022, 2023).

## 4. Conclusions

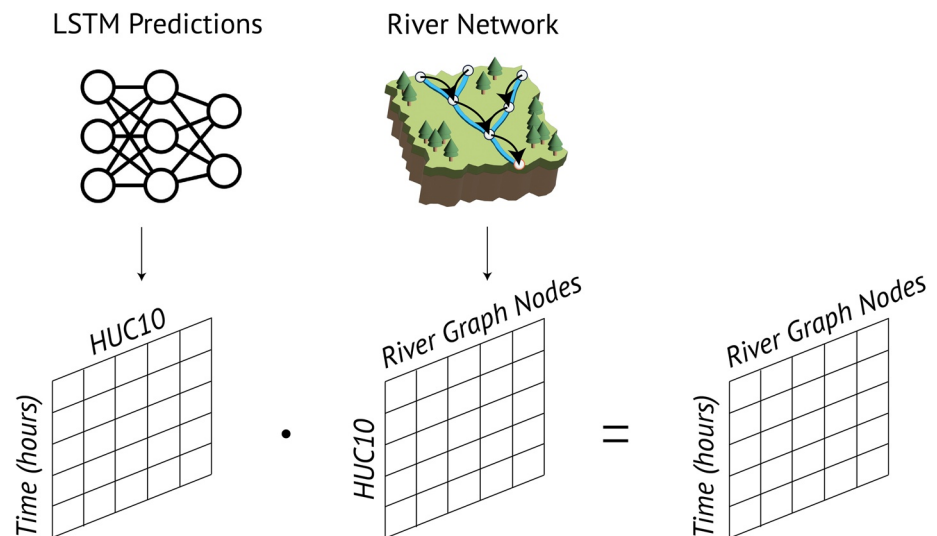
In this work, we used a combination of a pre-trained LSTM rainfall-runoff model and MC routing to create a learnable routing model, or, from the perspective of ML, a physics-informed GNN. This model predicts streamflow in stem rivers and learns river parameters throughout a river network, which is urgently needed to improve the next-generation large-scale hydrologic models. Because our framework is built on physical principles and

estimates widely used  $n$  values, it can be easily ported to work with other models. For example, the trained NN and the weights can be loaded into Fortran or C programs to support traditional hydrologic models or routing schemes, for example, (H. Li et al., 2013; Mizukami et al., 2016). Our synthetic experiments recovered the overall spatial pattern of  $n$  with moderate accuracy but could not recover the channel cross-sectional geometry parameter ( $q$ ). Furthermore, our synthetic experiments yielded promising results in recovering synthetic  $n$  and DA relationships, implying there is potential to learn reach-scale physics in the river network using differentiable modeling.

With real-world data, short-term training periods of downstream hydrographs can produce  $n$  parameterizations that improve long-term routing results, but may be insufficient to constrain the  $n$  values more precisely than a general spatial pattern. Eight weeks of real-world data produced decent long-term streamflow routing and improved upon approaches that did not use routing, yet training on different periods could result in somewhat different distributions. When looking at the  $n$  versus DA distribution attained by our trained model against USGS observations, we found that the  $n$  values agreed with the literature bounds for the area, but the absolute magnitudes may fluctuate depending on the training period and model mechanics. Besides using longer training periods to obtain  $n$  values that compromise across periods, future work should also consider whether  $n$  should be treated as dynamic in time. Further work can expand this analysis to other basins with different conditions (streams outside of the Ridge and Valley physiographic division of the CONUS) to see if the model can still identify their trends correctly. Reviewing the internal gage NSE scores over a full year of data showed a correlation between DA and the relative advantage of our routing scheme, highlighting the impacts of heterogeneity and flow convergence.

## Appendix A

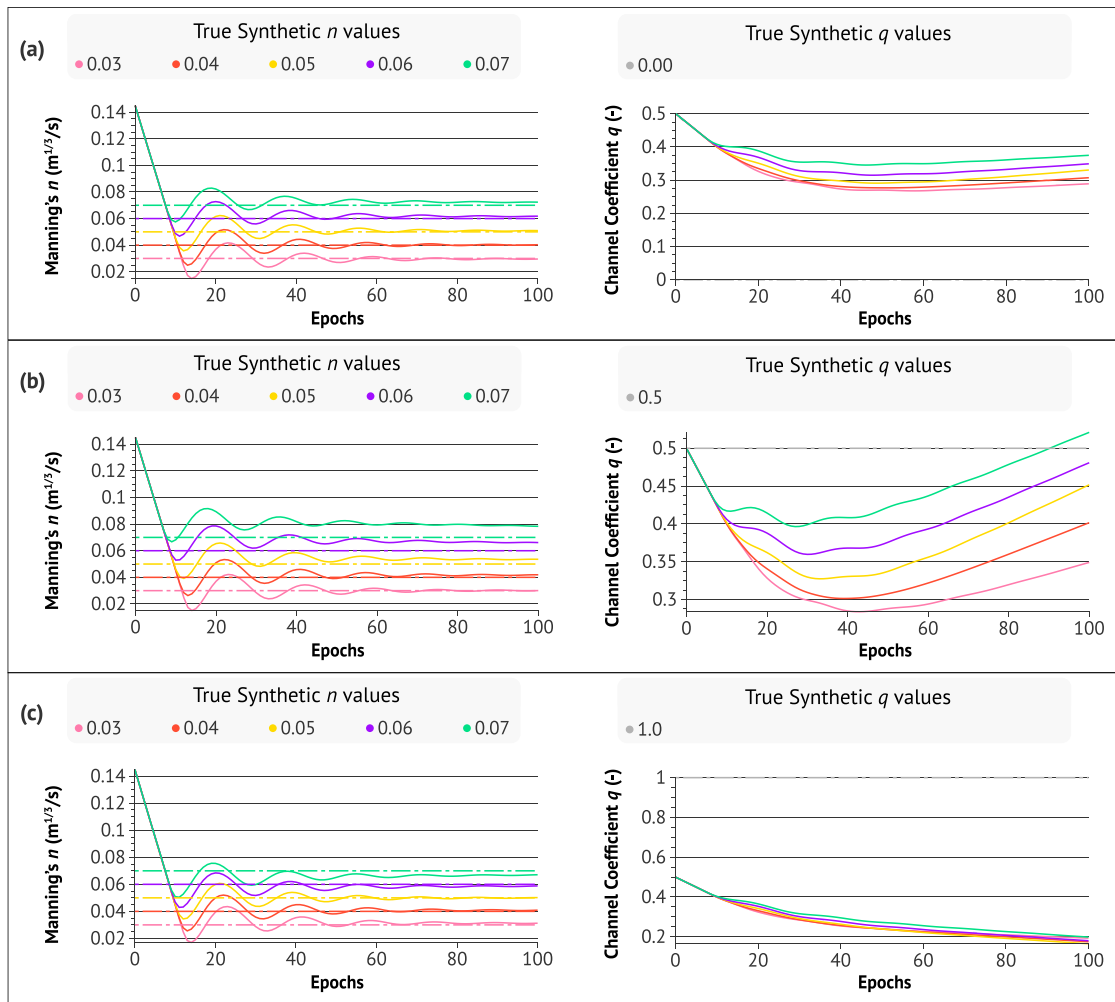
The contents of this appendix provide additional details on our methods. We include input attributes and meteorological forcings for the pretrained LSTM, and attributes for the Multilayer Perceptron network (MLP), with data sources. We also provide completed experimental results relevant to the Mean Integrated Gradient of each MLP parameter, and spatial  $q$  parameter recovery. Furthermore, we list the  $Q'$  NSE values, and Manning's  $n$  distributions, for different training periods.



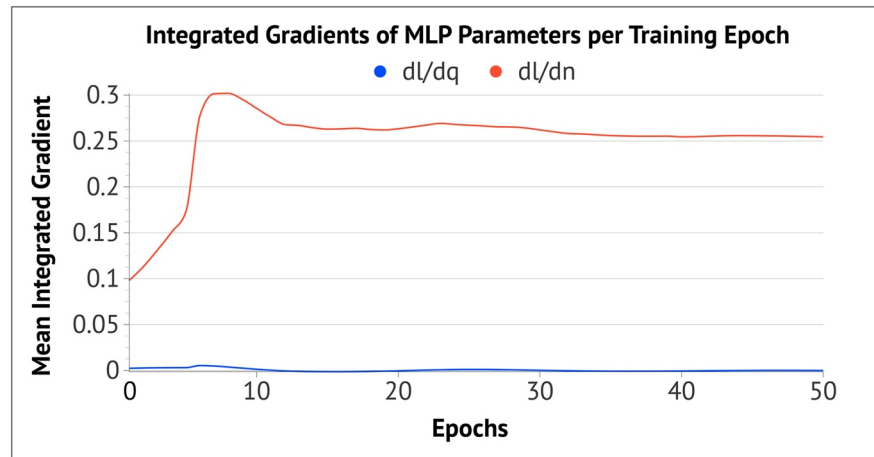
**Figure A1.** A visual representation of the mass-matrix multiplication to map long short-term memory HUC10 predictions into river segment hourly predictions.



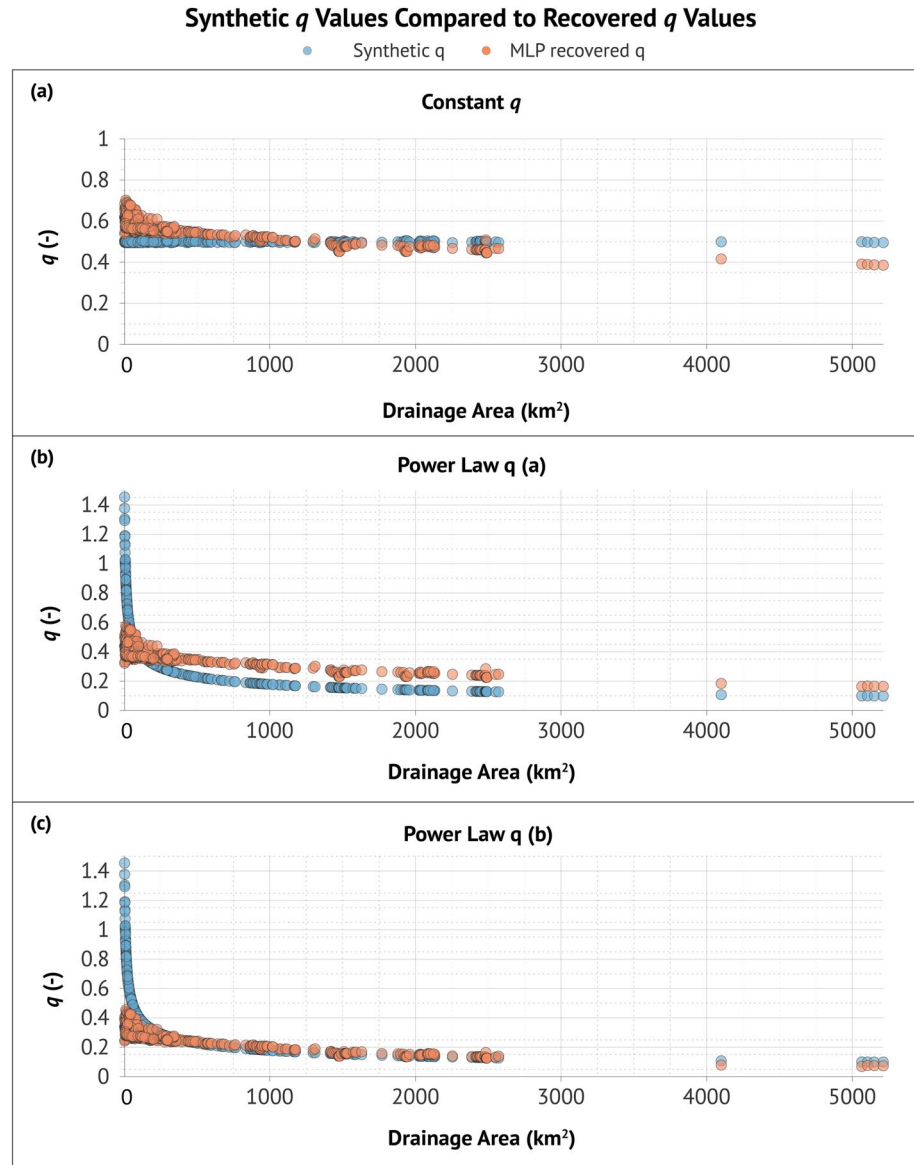
Single Parameter Recovery Experiments



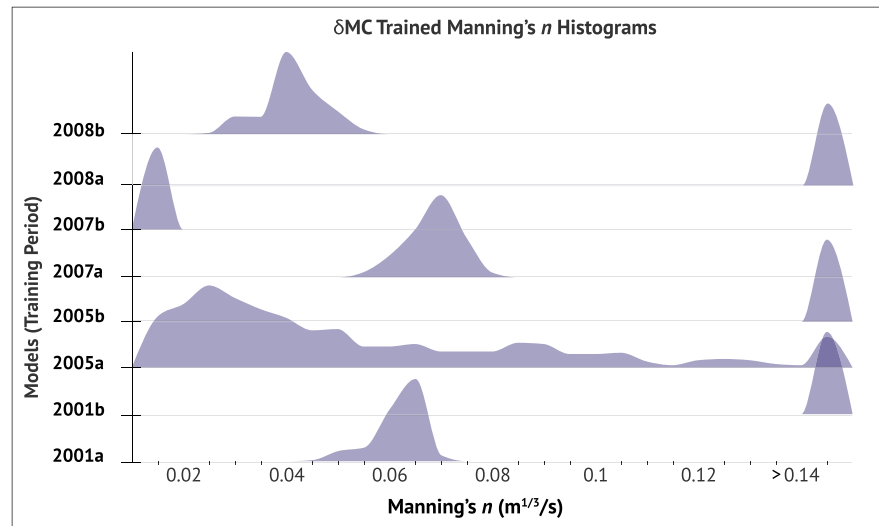
**Figure A2.** The synthetic parameter recovery of Manning's  $n$  after each epoch, with each colored line representing a different experiment with a unique “synthetic truth” combination of  $n$  and  $q$ . The figure is broken into (a), (b), and (c) to organize the experiments by the synthetic truth values used for  $q$  shown in the plot legend.



**Figure A3.** A line plot showing the Mean Integrated Gradient of each Multilayer Perceptron parameter ( $dL/dn$  and  $dL/dq$ ) per each epoch of model training.  $dL/dn$  and  $dL/dq$  were calculated from gradients accumulated when training  $\delta$ MC during period 2001a. A higher gradient value means the variable has more influence on parameter tuning.



**Figure A4.** Results from  $q$  parameter recovery experiments outlined in Section 2.6.1. We tried to recover both constant and distributed parameters but were unable to consistently recover the synthetic truth values.



**Figure A5.** Histograms visualizing the frequency of Manning's  $n$  values for all river reaches (582 total) for all eight graph neural network models. The lower bound is 0.01, while the upper bound contains all Manning's  $n$  values  $>0.14$ . The models often found extreme values for  $n$  when there was poor Nash-Sutcliffe Efficiency performance for their training periods.

**Table A1**

*The Attributes and Forcings Used by the Pre-Trained Long Short-Term Memory Rainfall-Runoff Model Used to Predict Streamflow (Links to the Data Can Be Found Below the Table)*

Attribute/Meteorological forcing	Unit	Data set	Citation
Mean elevation	m	SRTMGL1	Carabajal and Harding (2006)
Mean slope	Unitless	SRTMGL1	Carabajal and Harding (2006)
Basin area	km <sup>2</sup>	SRTMGL1	Carabajal and Harding (2006)
Dominant land cover	Class	MODIS	Friedl and Sulla-Menashe (2019)
Dominant land cover fraction	Percent	MODIS	Friedl and Sulla-Menashe (2019)
Forest fraction	Percent	MODIS	Friedl and Sulla-Menashe (2019)
Root depth (50)	m	MODIS	Friedl and Sulla-Menashe (2019)
Soil depth	m	MODIS	Friedl and Sulla-Menashe (2019)
Ksat (0–5)	log <sub>10</sub> (cm/hr)	POLARIS	Chaney et al. (2019)
Ksat (5–15)	log <sub>10</sub> (cm/hr)	POLARIS	Chaney et al. (2019)
Theta $s$ (0–5)	m <sup>3</sup> /m <sup>3</sup>	POLARIS	Chaney et al. (2019)
Theta $s$ (5–15)	m <sup>3</sup> /m <sup>3</sup>	POLARIS	Chaney et al. (2019)
Theta $r$ (5–15)	m <sup>3</sup> /m <sup>3</sup>	POLARIS	Chaney et al. (2019)
Ksat average (0–15)	log <sub>10</sub> (cm/hr)	POLARIS	Chaney et al. (2019)
Ksat $e$ (0–5)	cm/hr	POLARIS	Chaney et al. (2019)
Ksat $e$ (5–15)	cm/hr	POLARIS	Chaney et al. (2019)
Ksat average $e$ (0–15)	cm/hr	POLARIS	Chaney et al. (2019)
Theta average $s$ (0–15)	e <sup>m3/m3</sup>	POLARIS	Chaney et al. (2019)
Theta average $r$ (0–15)	e <sup>m3/m3</sup>	POLARIS	Chaney et al. (2019)
Porosity	Percent	GLHYMPS	Huscroft et al. (2018)
Permeability permafrost	m <sup>2</sup>	GLHYMPS	Huscroft et al. (2018)
Permeability permafrost (Raw)	m <sup>2</sup>	GLHYMPS	Huscroft et al. (2018)
Major number of dams	Unitless	GAGES-II	Falcone (2011)

**Table A1**

*Continued*

Attribute/Meteorological forcing	Unit	Data set	Citation
General purpose of dam	Unitless	National Inventory of Dams (NID)	US Army Corps of Engineers (2018)
Max of normal storage	Acre-ft	National Inventory of Dams (NID)	US Army Corps of Engineers (2018)
Standard deviation of normal storage	Unitless	National Inventory of Dams (NID)	US Army Corps of Engineers (2018)
Number of dams within river (2009)	Unitless	GAGES-II	Falcone (2011)
Normal storage (2009)	Acre-ft	National Inventory of Dams (NID)	US Army Corps of Engineers (2018)
Precipitation hourly total	kg/m <sup>2</sup>	NLDAS2	Xia et al. (2012)
Surface downward longwave radiation	W/m <sup>2</sup>	NLDAS2	Xia et al. (2012)
Surface downward shortwave radiation	W/m <sup>2</sup>	NLDAS2	Xia et al. (2012)
Pressure	Pa	NLDAS2	Xia et al. (2012)
Air temperature	K	NLDAS2	Xia et al. (2012)

*Note.* SRTMGL1: <https://doi.org/10.14358/PERS.72.3.287>, MODIS: <https://modis.gsfc.nasa.gov/data/dataproduct/mod12.php>, POLARIS: <https://doi.org/10.1029/2018WR022797>, GLHYMPS: <https://doi.org/10.5683/SP2/DLGXYO>, NID: <https://nid.usace.army.mil/>, NLDAS2: <https://ldas.gsfc.nasa.gov/nldas/v2/forcing>.

**Table A2**

*The Constant Attributes (A) Used by the Multilayer Perceptron to Predict  $n$  and  $q$ :  $n, q = NN(A)$*

Attribute	Unit
Reach Width	m
Average-Reach Elevation	m
Slope	m/m
Reach area	km <sup>2</sup>
Total drainage area	km <sup>2</sup>
Reach length	m
Sinuosity	m/m
Bank elevation	m

**Table A3**

*The  $\Sigma Q'$  ( $\tau = 9$ ) Nash-Sutcliffe Efficiency (NSE) Scores for All Eight Training Time Periods for the Furthest Downstream Gage*

	Periods							
	2001a	2001b	2005a	2005b	2007a	2007b	2008a	2008b
NSE	0.5956	0.3538	−0.7929	−0.1697	0.6835	0.0568	−0.4324	0.3798

*Note.* Since  $Q'$  routing is a pure forward simulation using the trained LSTM, we report the NSE values for each period.

## Conflict of Interest

This interest has been reviewed by the University in accordance with its Individual Conflict of Interest policy, for the purpose of maintaining the objectivity and the integrity of research at The Pennsylvania State University.



## Data Availability Statement

The data used in this study can be accessed via Zenodo (Bindas et al., 2023) while the differentiable routing model can be downloaded at Zenodo (Bindas & Shen, 2023). The LSTM streamflow model code (Feng et al., 2020; Ouyang et al., 2021) relevant to this work can be accessed via Zenodo (Shen, Fang, et al., 2021). For convenience, access to all of the above code is also compiled at the HydroDL research code hub (Shen, Bindas, et al., 2023). All data sets used are publicly available, including the GAGES-II data set (Falcone, 2011), NHDPlus (HorizonSystems, 2016), and NLDAS (Xia et al., 2012). Other data sources can be found in Table A1. Configuration files, and logging functionality, were built using Hydra (Yadan, 2019).

## Acknowledgments

We greatly appreciate the comments from editors and five anonymous reviewers whose comments have helped to improve this manuscript. The work was mainly supported by the Cooperative Institute for Research to Operations in Hydrology (CIROH) through the NOAA Cooperative Agreement with The University of Alabama (NA22NWS4320003), under subaward A22-0307-S003. Bindas would also like to acknowledge Penn State College of Engineering for providing University Graduate Fellowship during the first year of his graduate study. WP, JL and DF were supported by U.S. Department of Energy, Office of Science under Award DE-SC0016605. FR was supported by U.S. Department of Interior under award G21AC10563-00. Computing was partially supported by U.S. National Science Foundation Award PHY #2018280. KL and CS have financial interests in HydroSapient, Inc., a company which could potentially benefit from the results of this research.

## References

- Aboelyazeed, D., Xu, C., Hoffman, F. M., Liu, J., Jones, A. W., Rackauckas, C., et al. (2023). A differentiable, physics-informed ecosystem modeling and learning framework for large-scale inverse problems: Demonstration with photosynthesis simulations. *Biogeosciences*, 20(13), 2671–2692. <https://doi.org/10.5194/bg-20-2671-2023>
- Adnan, R. M., Petroselli, A., Heddad, S., Santos, C. A. G., & Kisi, O. (2021). Comparison of different methodologies for rainfall–runoff modeling: Machine learning vs conceptual approach. *Natural Hazards*, 105(3), 2987–3011. <https://doi.org/10.1007/s11069-020-04438-2>
- Arcement, G. J., & Schneider, V. R. (1989). *Guide for selecting Manning's roughness coefficients for natural channels and flood plains*. U.S. Geological Survey. (Water-Supply Paper No. 2339) Retrieved from <https://pubs.usgs.gov/wsp/2339/report.pdf>
- Baydin, A. G., Pearlmutter, B. A., Radul, A. A., & Siskind, J. M. (2018). Automatic differentiation in machine learning: A survey. *Journal of Machine Learning Research*, 18(153), 1–43.
- Beven, K. (2006). A manifesto for the equifinality thesis. *Journal of Hydrology*, 320(1–2), 18–36. <https://doi.org/10.1016/j.jhydrol.2005.07.007>
- Bindas, T., & Shen, C. (2023). mhp/dMC-Juniata-hydroDL2: v0.1.1. Version v0.1.1. Zenodo. <https://doi.org/10.5281/zenodo.10183449>
- Bindas, T., Tsai, W.-P., Liu, J., Rahmani, F., Feng, D., Bian, Y., et al. (2023). dMC-Juniata-hydroDL River Graph Dataset [Dataset]. Zenodo. Retrieved from <https://zenodo.org/records/10183429>
- Camporese, M., Paniconi, C., Putti, M., & Orlandini, S. (2010). Surface-subsurface flow modeling with path-based runoff routing, boundary condition-based coupling, and assimilation of multisource observation data. *Water Resources Research*, 46(2), W02512. <https://doi.org/10.1029/2008wr007536>
- Candela, A., Noto, L. V., & Aronica, G. (2005). Influence of surface roughness in hydrological response of semiarid catchments. *Journal of Hydrology*, 313(3), 119–131. <https://doi.org/10.1016/j.jhydrol.2005.01.023>
- Carabai, C. C., & Harding, D. J. (2006). SRTM C-Band and ICESat laser altimetry elevation comparisons as a function of tree cover and relief. *Photogrammetric Engineering and Remote Sensing*, 72(3), 287–298. <https://doi.org/10.14358/PERS.72.3.287>
- Chaney, N. W., Minasny, B., Herman, J. D., Nauman, T. W., Brungard, C. W., Morgan, C. L. S., et al. (2019). POLARIS soil properties: 30-m probabilistic maps of soil properties over the contiguous United States. *Water Resources Research*, 55(4), 2916–2938. <https://doi.org/10.1029/2018WR022797>
- Cunge, J. A. (1969). On the subject of a flood propagation computation method (Muskingum method). *Journal of Hydraulic Research*, 7(2), 205–230. <https://doi.org/10.1080/00221686909500264>
- David, C. H., Habets, F., Maidment, D. R., & Yang, Z.-L. (2011). RAPID applied to the SIM-France model. *Hydrological Processes*, 25(22), 3412–3425. <https://doi.org/10.1002/hyp.8070>
- David, C. H., Maidment, D. R., Niu, G.-Y., Yang, Z.-L., Habets, F., & Eijkhout, V. (2011). River network routing on the NHDPlus dataset. *Journal of Hydrometeorology*, 12(5), 913–934. <https://doi.org/10.1175/2011JHM1345.1>
- David, C. H., Yang, Z.-L., & Hong, S. (2013). Regional-scale river flow modeling using off-the-shelf runoff products, thousands of mapped rivers and hundreds of stream flow gauges. *Environmental Modelling and Software*, 42, 116–132. <https://doi.org/10.1016/j.envsoft.2012.12.011>
- Dottori, F., Szewczyk, W., Ciscar, J.-C., Zhao, F., Alfieri, L., Hirabayashi, Y., et al. (2018). Increased human and economic losses from river flooding with anthropogenic warming. *Nature Climate Change*, 8(9), 781–786. <https://doi.org/10.1038/s41558-018-0257-z>
- Douben, K.-J. (2006). Characteristics of river floods and flooding: A global overview. *Irrigation and Drainage*, 55(S1), S9–S21. <https://doi.org/10.1002/ird.239>
- Duan, S., Ullrich, P., & Shu, L. (2020). Using convolutional neural networks for streamflow projection in California. *Frontiers in Water*, 2. <https://doi.org/10.3389/frwa.2020.00028>
- Falcone, J. A. (2011). GAGES-II: Geospatial attributes of gages for evaluating streamflow [Dataset]. In *GAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow*. USGS Unnumbered Series. U.S. Geological Survey. <https://doi.org/10.3133/70046617>
- Fan, W., Ma, Y., Li, Q., He, Y., Zhao, E., Tang, J., & Yin, D. (2019). Graph neural networks for social recommendation. arXiv. <https://doi.org/10.48550/arXiv.1902.07243>
- Fang, K., Pan, M., & Shen, C. (2019). The value of SMAP for long-term soil moisture estimation with the help of deep learning. *IEEE Transactions on Geoscience and Remote Sensing*, 57(4), 2221–2233. <https://doi.org/10.1109/TGRS.2018.2872131>
- Fang, K., & Shen, C. (2017). Full-flow-regime storage-streamflow correlation patterns provide insights into hydrologic functioning over the continental US. *Water Resources Research*, 53(9), 8064–8083. <https://doi.org/10.1002/2016WR020283>
- Fang, K., Shen, C., Kifer, D., & Yang, X. (2017). Prolongation of SMAP to spatiotemporally seamless coverage of continental U.S. using a deep learning neural network. *Geophysical Research Letters*, 44(21), 11030–11039. <https://doi.org/10.1002/2017gl075619>
- Feng, D., Beck, H., Lawson, K., & Shen, C. (2023). The suitability of differentiable, physics-informed machine learning hydrologic models for ungauged regions and climate change impact assessment. *Hydrology and Earth System Sciences*, 27(12), 2357–2373. <https://doi.org/10.5194/hess-27-2357-2023>
- Feng, D., Fang, K., & Shen, C. (2020). Enhancing streamflow forecast and extracting insights using long-short term memory networks with data integration at continental scales. *Water Resources Research*, 56(9), e2019WR026793. <https://doi.org/10.1029/2019WR026793>
- Feng, D., Liu, J., Lawson, K., & Shen, C. (2022). Differentiable, learnable, regionalized process-based models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy. *Water Resources Research*, 58(10), e2022WR032404. <https://doi.org/10.1029/2022WR032404>
- France-Presse, A. (2022). At least 59 dead and millions stranded as floods devastate India and Bangladesh. The Guardian. Retrieved from <https://www.theguardian.com/world/2022/jun/18/at-least-18-dead-and-millions-stranded-as-floods-devastate-india-and-bangladesh>

- François, B., Schlef, K. E., Wi, S., & Brown, C. M. (2019). Design considerations for riverine floods in a changing climate—A review. *Journal of Hydrology*, 574, 557–573. <https://doi.org/10.1016/j.jhydrol.2019.04.068>
- Friedl, M., & Sulla-Menasse, D. (2019). MCD12Q1 MODIS/Terra+Aqua land cover type yearly L3 global 500 m SIN grid V006 [Dataset]. <https://doi.org/10.5067/MODIS/MCD12Q1.006>
- Getirana, A. C. V., Boone, A., Yamazaki, D., Decharme, B., Papa, F., & Mognard, N. (2012). The hydrological modeling and analysis platform (HyMAP): Evaluation in the Amazon basin. *Journal of Hydrometeorology*, 13(6), 1641–1665. <https://doi.org/10.1175/JHM-D-12-021.1>
- Gleason, C. J. (2015). Hydraulic geometry of natural rivers: A review and future directions. *Progress in Physical Geography*, 39(3), 337–360. <https://doi.org/10.1177/0309133314567584>
- Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (pp. 249–256). JMLR Workshop and Conference Proceedings. Retrieved from <https://proceedings.mlr.press/v9/glorot10a.html>
- Greuell, W., Andersson, J. C. M., Donnelly, C., Feyen, L., Gerten, D., Ludwig, F., et al. (2015). Evaluation of five hydrological models across Europe and their suitability for making projections under climate change. *Hydrology and Earth System Sciences Discussions*, 12(10), 10289–10330. <https://doi.org/10.5194/hessd-12-10289-2015>
- Harmen, G. (1955). The concept “hydraulic radius” in porous media. *Transactions of the AIME*, 204(01), 274–277. <https://doi.org/10.2118/479-G>
- He, M., Wu, S., Huang, B., Kang, C., & Gui, F. (2022). Prediction of total nitrogen and phosphorus in surface water by deep learning methods based on multi-scale feature extraction. *Water*, 14(10), 1643. <https://doi.org/10.3390/w14101643>
- Hochreiter, S. (1998). The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 06(02), 107–116. <https://doi.org/10.1142/S0218488598000094>
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- HorizonSystems. (2016). NHDPlus version 2 [Dataset]. Retrieved from [http://www.horizon-systems.com/nhdplus/NHDplusV2\\_home.php](http://www.horizon-systems.com/nhdplus/NHDplusV2_home.php)
- Hrnjica, B., Mehr, A. D., Jakupović, E., Crnković, A., & Hasanagić, R. (2021). Application of deep learning neural networks for nitrate prediction in the Klokot River, Bosnia and Herzegovina. In *2021 7th International Conference on Control* (pp. 1–6). Instrumentation and Automation (ICCIA). <https://doi.org/10.1109/ICCIA52082.2021.9403565>
- Huscroft, J., Gleason, T., Hartmann, J., & Börker, J. (2018). Compiling and mapping global permeability of the unconsolidated and consolidated Earth: GLobal HYdrogeology MaPS 2.0 (GLHYMPS 2.0). *Geophysical Research Letters*, 45(4), 1897–1904. <https://doi.org/10.1002/2017GL075860>
- International Panel on Climate Change (IPCC). (2012). Managing the risks of extreme events and disasters to advance climate change adaptation (p. 582). Retrieved from <https://www.ipcc.ch/report/managing-the-risks-of-extreme-events-and-disasters-to-advance-climate-change-adaptation/>
- Ji, X., Lesack, L., Melack, J. M., Wang, S., Riley, W. J., & Shen, C. (2019). Seasonal and inter-annual patterns and controls of hydrological fluxes in an Amazon floodplain lake with a surface-subsurface processes model. *Water Resources Research*, 55(4), 3056–3075. <https://doi.org/10.1029/2018WR023897>
- Jia, X., Zwart, J., Sadler, J., Appling, A., Oliver, S., Markstrom, S., et al. (2021). Physics-guided recurrent graph model for predicting flow and temperature in river networks. In *Proceedings of the 2021 SIAM International Conference on Data Mining (SDM)* (pp. 612–620). Society for Industrial and Applied Mathematics. <https://doi.org/10.1137/1.9781611976700.69>
- Johnson, J. M., Munasinghe, D., Eyelade, D., & Cohen, S. (2019). An integrated evaluation of the national water model (NWM)—Height above nearest drainage (HAND) flood mapping methodology. *Natural Hazards and Earth System Sciences*, 19(11), 2405–2420. <https://doi.org/10.5194/nhess-19-2405-2019>
- Kalyanapu, A. J., Burian, S. J., & McPherson, T. N. (2009). Effect of land use-based surface roughness on hydrologic model output. *Journal of Spatial Hydrology*, 9(2), 51–71.
- Khorashadi Zadeh, F., Nossent, J., Sarrazin, F., Pianosi, F., van Griensven, A., Wagener, T., & Bauwens, W. (2017). Comparison of variance-based and moment-independent global sensitivity analysis approaches by application to the SWAT model. *Environmental Modelling and Software*, 91, 210–222. <https://doi.org/10.1016/j.envsoft.2017.02.001>
- Khoshkalam, Y., Rousseau, A. N., Rahmani, F., Shen, C., & Abbasnezhadi, K. (2023). Applying transfer learning techniques to enhance the accuracy of streamflow prediction produced by long Short-term memory networks with data integration. *Journal of Hydrology*, 622, 129682. <https://doi.org/10.1016/j.jhydrol.2023.129682>
- Kingma, D. P., & Ba, J. (2017). Adam: A method for stochastic optimization. arXiv. <https://doi.org/10.48550/arXiv.1412.6980>
- Koks, E. E., & Thissen, M. (2016). A multiregional impact assessment model for disaster analysis. *Economic Systems Research*, 28(4), 429–449. <https://doi.org/10.1080/09535314.2016.1232701>
- Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019). Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets. *Hydrology and Earth System Sciences*, 23(12), 5089–5110. <https://doi.org/10.5194/hess-23-5089-2019>
- Leopold, L. B., & Maddock, T., Jr. (1953). *The hydraulic geometry of stream channels and some physiographic implications* (p. 252). USGS Professional Paper. <https://doi.org/10.3133/pp252>
- Leshno, M., Lin, V. Y., Pinkus, A., & Schocken, S. (1993). Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6), 861–867. [https://doi.org/10.1016/S0893-6080\(05\)80131-5](https://doi.org/10.1016/S0893-6080(05)80131-5)
- Li, H., Wigmosta, M. S., Wu, H., Huang, M., Ke, Y., Coleman, A. M., & Leung, L. R. (2013). A physically based runoff routing model for land surface and Earth system models. *Journal of Hydrometeorology*, 14(3), 808–828. <https://doi.org/10.1175/JHM-D-12-015.1>
- Li, H.-Y., Tan, Z., Ma, H., Zhu, Z., Abeshu, G. W., Zhu, S., et al. (2022). A new large-scale suspended sediment model and its application over the United States. *Hydrology and Earth System Sciences*, 26(3), 665–688. <https://doi.org/10.5194/hess-26-665-2022>
- Lin, G.-Y., Chen, H.-W., Chen, B.-J., & Yang, Y.-C. (2022). Characterization of temporal PM2.5, nitrate, and sulfate using deep learning techniques. *Atmospheric Pollution Research*, 13(1), 101260. <https://doi.org/10.1016/j.apr.2021.101260>
- Liu, J., Hughes, D., Rahmani, F., Lawson, K., & Shen, C. (2023). Evaluating a global soil moisture dataset from a multitask model (GSM3 v1.0) with potential applications for crop threats. *Geoscientific Model Development*, 16(5), 1553–1567. <https://doi.org/10.5194/gmd-16-1553-2023>
- Liu, J., Rahmani, F., Lawson, K., & Shen, C. (2022). A multiscale deep learning model for soil moisture integrating satellite and in situ data. *Geophysical Research Letters*, 49(7), e2021GL096847. <https://doi.org/10.1029/2021GL096847>
- Liu, L., Ao, T., Zhou, L., Takeuchi, K., Gusyev, M., Zhang, X., et al. (2022). Comprehensive evaluation of parameter importance and optimization based on the integrated sensitivity analysis system: A case study of the BTOP model in the upper Min River Basin, China. *Journal of Hydrology*, 610, 127819. <https://doi.org/10.1016/j.jhydrol.2022.127819>
- Mangukiyi, N. K., Sharma, A., & Shen, C. (2023). How to enhance hydrological predictions in hydrologically distinct watersheds of the Indian subcontinent? *Hydrological Processes*, 37(7), e14936. <https://doi.org/10.1002/hyp.14936>

- Martinez, G. F., & Gupta, H. V. (2010). Toward improved identification of hydrological models: A diagnostic evaluation of the “abcd” monthly water balance model for the conterminous United States. *Water Resources Research*, 46(8), W08507. <https://doi.org/10.1029/2009WR008294>
- Mays, L. W. (2010). *Water resources engineering* (2nd ed.). Wiley.
- Mays, L. W. (2019). *Water resources engineering* (3rd ed.). Wiley. Retrieved from <https://www.wiley.com/en-us/Water+Resources+Engineering+%2C+3rd+Edition-p-9781119493167>
- Meyal, A. Y., Versteeg, R., Alper, E., Johnson, D., Rodzianko, A., Franklin, M., & Wainwright, H. (2020). Automated cloud based long short-term memory neural network based SWE prediction. *Frontiers in Water*, 2. <https://doi.org/10.3389/frwa.2020.574917>
- Mizukami, N., Clark, M. P., Sampson, K., Nijssen, B., Mao, Y., McMillan, H., et al. (2016). mizuRoute version 1: A river network routing tool for a continental domain water resources applications. *Geoscientific Model Development*, 9(6), 2223–2238. <https://doi.org/10.5194/gmd-9-2223-2016>
- Moore, R. B., & Dewald, T. G. (2016). The road to NHDPlus—Advancements in digital stream networks and associated catchments. *JAWRA Journal of the American Water Resources Association*, 52(4), 890–900. <https://doi.org/10.1111/1752-1688.12389>
- Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I—A discussion of principles. *Journal of Hydrology*, 10(3), 282–290. [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)
- Natural Environment Research Council. (1975). In *Flood routing studies* (Vol. 3). Wallingford, England: Institute of Hydrology.
- Orlandini, S. (2002). On the spatial variation of resistance to flow in upland channel networks. *Water Resources Research*, 38(10), 15–15-14. <https://doi.org/10.1029/2001wr001187>
- Orlandini, S., & Rosso, R. (1998). Parameterization of stream channel geometry in the distributed modeling of catchment dynamics. *Water Resources Research*, 34(8), 1971–1985. <https://doi.org/10.1029/98wr00257>
- O, S., & Orth, R. (2021). Global soil moisture data derived through machine learning trained with in-situ measurements. *Scientific Data*, 8(1), 170. <https://doi.org/10.1038/s41597-021-00964-1>
- Ouyang, W., Lawson, K., Feng, D., Ye, L., Zhang, C., & Shen, C. (2021). Continental-scale streamflow modeling of basins with reservoirs: Towards a coherent deep-learning-based strategy. *Journal of Hydrology*, 599, 126455. <https://doi.org/10.1016/j.jhydrol.2021.126455>
- Pan, M., & Wood, E. F. (2013). Inverse streamflow routing. *Hydrology and Earth System Sciences*, 17(11), 4577–4588. <https://doi.org/10.5194/hess-17-4577-2013>
- Ponce, V. M. (1986). Diffusion wave modeling of catchment dynamics. *Journal of Hydraulic Engineering*, 112(8), 716–727. [https://doi.org/10.1061/\(ASCE\)0733-9429\(1986\)112:8\(716\)](https://doi.org/10.1061/(ASCE)0733-9429(1986)112:8(716))
- Prein, A. F., Rasmussen, R. M., Ikeda, K., Liu, C., Clark, M. P., & Holland, G. J. (2017). The future intensification of hourly precipitation extremes. *Nature Climate Change*, 7(1), 48–52. <https://doi.org/10.1038/nclimate3168>
- Rahmani, F., Appling, A., Feng, D., Lawson, K., & Shen, C. (2023). Identifying structural priors in a hybrid differentiable model for stream water temperature modeling. *Water Resources Research*, 59(12), e2023WR034420. <https://doi.org/10.1029/2023WR034420>
- Rahmani, F., Lawson, K., Ouyang, W., Appling, A., Oliver, S., & Shen, C. (2021). Exploring the exceptional performance of a deep learning stream temperature model and the value of streamflow data. *Environmental Research Letters*, 16(2), 024025. <https://doi.org/10.1088/1748-9326/abd501>
- Rahmani, F., Shen, C., Oliver, S., Lawson, K., & Appling, A. (2021). Deep learning approaches for improving prediction of daily stream temperature in data-scarce, unmonitored, and dammed basins. *Hydrological Processes*, 35(11), e14400. <https://doi.org/10.1002/hyp.14400>
- Regan, R. S., Markstrom, S. L., Hay, L. E., Viger, R. J., Norton, P. A., Driscoll, J. M., & LaFontaine, J. H. (2018). Description of the national hydrologic model for use with the precipitation-runoff modeling system (PRMS) (No. 6-B9). In *Techniques and methods*. U.S. Geological Survey. <https://doi.org/10.3133/tm6B9>
- Rice, D. (2019). *Mississippi river flood is longest-lasting in over 90 years, since “Great Flood” of 1927*. USA Today. Retrieved from <https://www.usatoday.com/story/news/nation/2019/05/28/mississippi-river-flooding-longest-lasting-since-great-flood-1927/1261049001/>
- Saha, G. K., Rahmani, F., Shen, C., Li, L., & Cibir, R. (2023). A deep learning-based novel approach to generate continuous daily stream nitrate concentration for nitrate data-sparse watersheds. *Science of the Total Environment*, 878, 162930. <https://doi.org/10.1016/j.scitotenv.2023.162930>
- Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., et al. (2023). Differentiable modelling to unify machine learning and physical models for geosciences. *Nature Reviews Earth and Environment*, 4(8), 552–567. <https://doi.org/10.1038/s43017-023-00450-9>
- Shen, C., Bindas, T., Song, Y., Feng, D., & Sawadekar, K. (2023). HydroDL Docs. Retrieved from <https://mhpi.github.io/>
- Shen, C., Chen, X., & Laloy, E. (2021). Editorial: Broadening the use of machine learning in hydrology. *Frontiers in Water*, 3. <https://doi.org/10.3389/frwa.2021.681023>
- Shen, C., Fang, K., Feng, D., & Bindas, T. (2021). mhpi/hydroDL: MHPI-hydroDL [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.5015120>
- Shen, C., & Lawson, K. (2021). Applications of deep learning in hydrology. In *Deep learning for the Earth Sciences* (pp. 283–297). John Wiley & Sons, Ltd. <https://doi.org/10.1002/9781119646181.ch19>
- Shen, C., Niu, J., & Fang, K. (2014). Quantifying the effects of data integration algorithms on the outcomes of a subsurface–land surface processes model. *Environmental Modelling and Software*, 59, 146–161. <https://doi.org/10.1016/j.envsoft.2014.05.006>
- Shen, C., Niu, J., & Phanikumar, M. S. (2013). Evaluating controls on coupled hydrologic and vegetation dynamics in a humid continental climate watershed using a subsurface–Land surface processes model. *Water Resources Research*, 49(5), 2552–2572. <https://doi.org/10.1002/wrcr.20189>
- Shen, C., & Phanikumar, M. S. (2010). A process-based, distributed hydrologic model based on a large-scale method for surface–subsurface coupling. *Advances in Water Resources*, 33(12), 1524–1541. <https://doi.org/10.1016/j.advwatres.2010.09.002>
- Shen, C., Riley, W. J., Smithgall, K. M., Melack, J. M., & Fang, K. (2016). The fan of influence of streams and channel feedbacks to simulated land surface water and carbon dynamics. *Water Resources Research*, 52(2), 880–902. <https://doi.org/10.1002/2015WR018086>
- Sun, A. Y., Jiang, P., Mudunuru, M. K., & Chen, X. (2021). Explore spatio-temporal learning of large sample hydrology using graph neural networks. *Water Resources Research*, 57(12), e2021WR030394. <https://doi.org/10.1029/2021WR030394>
- Sun, A. Y., Jiang, P., Yang, Z.-L., Xie, Y., & Chen, X. (2022). A graph neural network (GNN) approach to basin-scale river network learning: The role of physics-based connectivity and data fusion. *Hydrology and Earth System Sciences Discussions*, 26(19), 5163–5184. <https://doi.org/10.5194/hess-26-5163-2022>
- Tsai, W.-P., Feng, D., Pan, M., Beck, H., Lawson, K., Yang, Y., et al. (2021). From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling. *Nature Communications*, 12(1), 5988. <https://doi.org/10.1038/s41467-021-26107-z>
- US Army Corps of Engineers. (2018). National inventory of dams (NID) [Dataset]. Retrieved from <https://nid.sec.usace.army.mil/>
- USGS ScienceBase-Catalog. (2022). National elevation dataset (NED) [Dataset]. Retrieved from <https://www.sciencebase.gov/catalog/item/4fcf8fd4e4b0c7fe80e81504>
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., et al. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3), 261–272. <https://doi.org/10.1038/s41592-019-0686-2>

- Winsemius, H. C., Aerts, J. C. J. H., van Beek, L. P. H., Bierkens, M. F. P., Bouwman, A., Jongman, B., et al. (2016). Global drivers of future river flood risk. *Nature Climate Change*, 6(4), 381–385. <https://doi.org/10.1038/nclimate2893>
- Wunsch, A., Liesch, T., & Broda, S. (2022). Deep learning shows declining groundwater levels in Germany until 2100 due to climate change. *Nature Communications*, 13(1), 1221. <https://doi.org/10.1038/s41467-022-28770-2>
- Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., et al. (2009). NLDAS primary forcing data L4 hourly 0.125 x 0.125 degree V002 (NLDAS\_FORA0125\_H) [Dataset]. Goddard Earth Sciences Data and Information Services Center (GES DISC). <https://doi.org/10.5067/6J5LHHOZH4>
- Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., et al. (2012). Continental-scale water and energy flux analysis and validation for the North American land data assimilation system project phase 2 (NLDAS-2): 1. Intercomparison and application of model products. *Journal of Geophysical Research*, 117(D3), D03109. <https://doi.org/10.1029/2011JD016048>
- Xiang, Z., Yan, J., & Demir, I. (2020). A rainfall-runoff model with LSTM-based sequence-to-sequence learning. *Water Resources Research*, 56(1), e2019WR025326. <https://doi.org/10.1029/2019WR025326>
- Yadan, O. (2019). Hydra—A framework for elegantly configuring complex applications. Retrieved from <https://github.com/facebookresearch/hydra>
- Ye, A., Zhou, Z., You, J., Ma, F., & Duan, Q. (2018). Dynamic Manning's roughness coefficients for hydrological modelling in basins. *Hydrology Research*, 49(5), 1379–1395. <https://doi.org/10.2166/nh.2018.175>
- Zhi, W., Feng, D., Tsai, W.-P., Sterle, G., Harpold, A., Shen, C., & Li, L. (2021). From hydrometeorology to river water quality: Can a deep learning model predict dissolved oxygen at the continental scale? *Environmental Science and Technology*, 55(4), 2357–2368. <https://doi.org/10.1021/acs.est.0c06783>
- Zhu, F., Li, X., Qin, J., Yang, K., Cuo, L., Tang, W., & Shen, C. (2021). Integration of multisource data to estimate downward longwave radiation based on deep neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, 60, 1–15. <https://doi.org/10.1109/TGRS.2021.3094321>