#### Research papers

Continental-scale streamflow modeling of basins with reservoirs: towards a coherent deep-learning-based strategy

Wenyu Ouyang, Kathryn Lawson, Dapeng Feng, Lei Ye, Chi Zhang, Chaopeng Shen

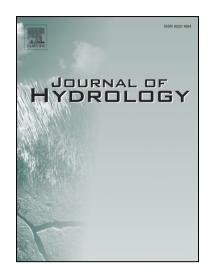
PII: S0022-1694(21)00502-3

DOI: https://doi.org/10.1016/j.jhydrol.2021.126455

Reference: HYDROL 126455

To appear in: Journal of Hydrology

Received Date: 4 January 2021 Revised Date: 10 May 2021 Accepted Date: 12 May 2021



Please cite this article as: Ouyang, W., Lawson, K., Feng, D., Ye, L., Zhang, C., Shen, C., Continental-scale streamflow modeling of basins with reservoirs: towards a coherent deep-learning-based strategy, *Journal of Hydrology* (2021), doi: https://doi.org/10.1016/j.jhydrol.2021.126455

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

© 2021 Elsevier B.V. All rights reserved.

#### 1 Continental-scale streamflow modeling of basins with reservoirs: towards a coherent

2 deep-learning-based strategy

3

4 Wenyu Ouyang<sup>1</sup>, Kathryn Lawson<sup>2</sup>, Dapeng Feng<sup>2</sup>, Lei Ye<sup>1</sup>, Chi Zhang<sup>1</sup>, Chaopeng Shen<sup>2,\*</sup>

5

- 6 <sup>1</sup> School of Hydraulic Engineering, Dalian University of Technology, Dalian, China
- 7 <sup>2</sup> Civil and Environmental Engineering, Pennsylvania State University, University Park, P.
- 8 USA

9

11

12

13 14 15

16

21 22 23

242526

27

28

29 30

31

10 Abstract

A large fraction of major waterways have dams influen ing greamflow, which must be streamflow prediction for accounted for in large-scale hydrologic modeling. However, basins with dams is challenging for various modeling approaches, especially at large scales. Here we examined which types of dammed basins of uld be well represented by long shortterm memory (LSTM) models using readil a labe information, and delineated the m 3557 basins (83% dammed) over the remaining challenges. We analyzed contiguous United States and noted stong it pact of reservoir purposes, degree of regulation (dor), and diversion on streamflow no eling. While a model trained on a widely-used reference-basin dataset performed poon, for non-reference basins, the model trained on the whole dataset presented a med an Nash-outcliffe efficiency coefficient (NSE) of 0.74. The zero-dor, small-dor (with storage c approximately a month of average streamflow or less), and large-dor basins were four to have distinct behaviors, so migrating models between categories yielded catastrophic results, which means we must not treat small-dor basins as reference ones. However, air with pooled data from different sets yielded open.

NSFs of 0.72, 0.79 and 0.6 for these respective groups, noticeably stronger than existing atrategy where smaller dams (storing models. These results support a coherent modeling strategy where smaller dams (storing about a month of average streamflow or less) are modeled implicitly as part of basin rainfallrunoff processes: being large-dor reservoirs of certain types can be represented explicitly. However, danged pasins must be present in the training dataset. Future work should examing rac modeling of large reservoirs for fire protection and irrigation, hydroelectric eneral on, and flood control.

32

<sup>\*</sup> corresponding author. Chaopeng Shen: cshen@engr.psu.edu

#### 1. Introduction.

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

34

Two-thirds of the longest rivers in the world are not flowing freely (Grill et al., 2019): more than 800,000 dammed reservoirs impede the world's rivers, including 90,000 in the United States (International Rivers, 2007). Dams exert significant control on streamflows by changing the magnitude and timing of the discharges (Gutenson et al., 2020). The ability to anticipate upstream reservoir operations at a daily scale has significant operational value for optimal water resources management.

For large-scale hydrologic modeling at the daily scale, we nee and tractable methods to account for the influence of small and large reserve is of eamflow. One may use a reservoir-centric modeling approach, in which each, eser pir leeds to be represented explicitly with its own characteristics, operational rules, to see, inflow, and outflow. This approach may not scale well to large scales, however as there may be dozens or even hundreds of reservoirs upstream of the et a large basin. A different approach would be basin-centric (or grid-centric, also called Imped), in which all the reservoirs in a subbasin (or a computational gridcell) are gruped together into one unit in the river routing module. Apparently, the basin-centri (o lumbed) paradigm can vastly reduce modeling complexity 2008). Alternatively, a mixed approach can be taken where (Ehsani et al., 2016; Pa some reservoirs at lum, and while some others are explicitly represented. Current large-scale hydrologic models such as the National Water Model (NWM) (Gochis et al., 2018), or land surface Tyon logic models with routing schemes, e.g. the Community Land Model (Lawrence et al. 20 has mulate some major reservoirs and make the habitual assumption of ignoring the reservoirs. The questions are then: (i) What kinds of reservoirs can be modeled in a lumped fashion and what kind cannot? (ii) Can we ignore the impacts of small reservoirs and assume they are behaviorally similar to undammed basins?

It has been difficult to reliably obtain strong model performance for dammed basins using a rule-based system at large scales. From a literature survey (see more details in Appendix Table S1), it seems difficult to obtain Nash-Sutcliffe model efficiency coefficient

(NSE) values that are higher than 0.65 by assuming generic reservoir operational schemes (Biemans et al., 2011; Hanasaki et al., 2006; Shin et al., 2019; Voisin et al., 2013). Hanasaki et al. (2006) derived a demand-driven approach for global reservoir routing and laid the foundation for subsequent developments, showing error reduction compared to no-reservoir simulations, but no NSE was reported. Voison et al. (2013) improved upon the formulation from Hanasaki et al. (2006) to the heavily dammed Columbia River Basin and reported decent correlation but mostly negative NSEs, indicating substantial biases. Unlike ieno rel ase schemes, empirically derived target storage-release functions can be para net rized for individual reservoirs with sufficiently long observational records of inflows, and storage levels, and can reproduce observed flows more accurately t al., 2020; Turner et al., 2020; Wu and Chen, 2012; Yassin et al., 2019; Zaja et al., **1**017; Zhao et al., 2016). Yassin et al. (2019) used piecewise-linear relationships between eservoir storage, inflow, and release to describe reservoir policies and obtain describe reservoirs and obtain describe reservoirs across the globe. Zajac et al. (2017) rep a paximum NSE of 0.61 for 390 stations around the world. Although these results replied a significant progress in research, further research was still needed to inform whether these improvements were robust when simulated inflows from the hydrologic models, ather than observed inflows, were used as the input to reservoir al., 2020). In addition, one can certainly argue the current modules at large scale performance levels left from for improvement, which can provide better utility for practical applications.

applied to Lablish data-driven rules that relate reservoir storage, inflow, and release data. Ehst size at al. (2016) used ANNs to predict daily release using previous days' reservoir storage volume along with inflow and release measurements, and reported an NSE of 0.86. Yang et al. (2019) similarly applied recurrent neural networks, using inflow and water storage as inputs, to simulate the daily operation of three multi-purpose reservoirs located in one basin, and reported an NSE value over 0.85. However, the use of recent storage and inflow data is akin to a form of data assimilation and is known to greatly improve simulations for short-term

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

forecast (Feng et al., 2020a), but we do not use recent observations here as our objective is long-term projection. In addition, the existing generally-available reservoir databases (Lehner et al., 2011; Mulligan et al., 2020; Patterson and Doyle, 2018) mainly provide information on dam design specifications or operational details for some of the most significant reservoirs, which is not available for large-scale modeling in dammed basins.

Recently, the long short-term memory (LSTM) network (Hochreiter and Schmidhuber, 1997), a deep learning (DL) algorithm, has been applied to explore the a streamflow in basins across the CONUS. It is relatively inexpensive (in terms on time) to apply at large spatial scales, and has grown to be a well-established hydrological g tool (Shen, 2018). LSTM-based models can effectively learn streamflow name and have shown superior performance compared to other hydrological beng max models (Ayzel et al., 2020; Feng et al., 2020a; Kratzert et al., 2019b). For example, Kratzert et al. (2019b) reported that the median NSE value in the evaluation period could each 0.74 for a 531-basin subset of the eol logy for Large-Sample Studies (CAMELS) 671-basin Catchment Attributes and dataset using the forcing data from No the American Land Data Assimilation (NLDAS) system. More recently, Feng et al. (2020) improved the forecast NSE median to 0.86 with the addition of a data integration kernel inicorporated recent discharge observations. However, the e studies were based on, is composed of basins that are CAMELS dataset, whi considered to be refer ace" or undisturbed basins, which have minimal anthropogenic impacts (i.e. mining land use changes, minimal human water withdrawals) (Addor et al., 2017; 15). To our knowledge, there is no systematic knowledge regarding how Newma ns in basins with significant human modifications such as reservoirs or water LST n, especially at large scales.

Here we followed a divide-and-conquer approach to tackle the difficult problem of long-term daily streamflow prediction from dammed basins, and to delineate where challenges reside. We addressed the following questions: (1) Given only generally-available reservoir information, how well can LSTM networks make long-term daily streamflow predictions for basins with reservoirs across the entire CONUS? (2) How differently do basins with or without

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

reservoirs of different sizes function in streamflow --- how much error are we making if we simply ignore small reservoirs and treat those basins with small reservoirs as reference basins? (3) What kinds of reservoirs (purpose, size, diversion) can be well modeled in a lumped fashion and what kinds cannot? These questions have not been answered in the literature and the answers will help the community to devise an informed and coherent modeling strategy. We further provide experiences to the community on how to best form an appropriate training dataset, e.g., whether we should include basins with or without reservoirs as whether we should stratify basins into different categories based on reservoir characteris its, or simply group them together.

#### 2. Methods

As an overview, LSTM-based models were strained to predict long-term daily streamflow from basins with or without receivors. The inputs include atmospheric forcing time series data and static basin attributes who lographic attributes and anthropogenic influences). We trained the models using wrious subsets from a newly compiled 3557-basin dataset across the CONUS as well as the CAMELS dataset. Basins with complete streamflow records from 1 January 1990 through 3. December 2009 were selected from the Geospatial Attributes of Gages for Evaluating streamflow II (GAGES-II) dataset (Falcone, 2011). Below we provide the details of the procedures.

#### 2.1. STM

Long Short-Term Memory (LSTM) networks are a special kind of recurrent neural network (RNN) which can both learn from sequential data and address the notorious exploding and/or vanishing gradient problem (Hochreiter, 1998). These networks are composed of memory cells, the keys to which are the "cell states" and "gates" that control information flow within the LSTM algorithm. Cell states allow information to be stored over long time periods,

which is important for modeling catchment processes like snow, subsurface flow, and reservoir storage. Based on the input of the current time step and the output from the previous one, a "forget gate" decides what information is going to be removed from the existing cell state. Next, a sigmoid layer and a tanh layer are applied as an "input gate" to update the cell state. Finally, the cell state is put through a tanh function and multiplied by the output of the sigmoid "output gate" to determine the final output.

There are different formulations of LSTM-based models. Kratzert et a an N-to-1 model to predict streamflow, which means that the input was a mu and the output was a one-step variable. An N-to-M LSTM-base also called a sequence-to-sequence model, was employed to predict multi-time-sta eamflows by Xiang et al. (2020). In the present study, following Feng et al. (2020), we trained a CONUS-scale N-to-N model using meteorological forcings and static attribute of the basins to predict daily discharge. Here we did not use discharge from previous days as inputs. We trained the model but for inference, we ran the model in a single on sequences of a fixed length (365 d forward pass through the full time period This procedure means that during training, the LSTM has no context for the initial input teps of each sequence. However, in our preliminary analysis, we added a warm-up period of a varatto not have any noticeable impact. Thus we neglected the warm-up period for perform ce reasons. The N-to-N model had significant advantages in efficiency, and could reach convergence for the 671 basins in the CAMELS dataset with 10 years of training data in 69 minutes on an NVIDIA 1080 Ti graphical processing unit (GPU). In r, the model was able to be trained on 10-year data for the entire 3557-basin dataset this pap ence was achieved (300 epochs) in 427 minutes of computational time. In our until randomly sampled for sites and training periods to form mini-batches and we defined the total number of iterations in an epoch as corresponding to the probability that 99% of the time periods of all basins are picked in the epoch.

The forward propagation equations of the present LSTM-based model can be summarized as the following (see Figure S1 in Appendix for more details), based on the notations in Fang et al. (2020).

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

174 
$$x^{(t)} = ReLU(W_{xx}x_0^{(t)} + b_{xx})$$
 (1)

175 
$$f^{(t)} = \sigma(D(W_{fx}x^{(t)}) + D(W_{fh}h^{(t-1)}) + b_f)$$
 (2)

176 
$$i^{(t)} = \sigma(D(W_{ix}x^{(t)}) + D(W_{ih}h^{(t-1)}) + b_i)$$
 (3)

177 
$$g^{(t)} = tanh(D(W_{ax}x^{(t)}) + D(W_{ah}h^{(t-1)}) + b_a)$$
 (4)

178 
$$o^{(t)} = \sigma(D(W_{ox}x^{(t)}) + D(W_{oh}h^{(t-1)}) + b_o)$$
 (5)

179 
$$s^{(t)} = f^{(t)} \odot s^{(t-1)} + i^{(t)} \odot g^{(t)}$$
 (6)

$$h^{(t)} = \tanh(s^{(t)}) \odot o^{(t)} \tag{7}$$

181 
$$y^{(t)} = W_{hy}h^{(t)} + b_y \tag{8}$$

where  $x_0^{(t)}$  is the vector of raw inputs for the time step t,  $x^{(t)}$  is the input vector in the LSTM cell, ReLU is the rectified linear unit,  $\sigma$  is the sigmoid activation function D is the dropout operator,  $\odot$  denotes pointwise multiplication, W's are network weights, b's are bias parameters,  $g^{(t)}$  is the output of the input node,  $f^{(t)}$ ,  $i^{(t)}$ , and  $s^{(t)}$  are respectively the forget, input, and output gates,  $s^{(t)}$  represents the states of memory cells,  $h^{(t)}$  represents hidden states, and  $y^{(t)}$  is the predicted output which is compared to streamflow observations.

The static catchment attributes were concarabled with the meteorological inputs at each time step to produce the input vector. To reduce overfitting, we employed dropout regularization, which stochastically sets same network connections to zero. Here, *D* applies dropout with constant dropout passes to recurrent connections, i.e., the connections that are set to zero stay the same throughout reach training instance. This kind of dropout over recurrent connections allows the network to be treated as a Bayesian network (Gal and Ghahramani, 2016). In addition, a nonlinear transformation with a linear function and rectified linear unit (ReLU) was accled in the first input layer, following Fang et al. (2020). This was used because without the injut transformation layer, some weights of inputs would be directly set to 0 after dropout and lead to information loss. The network outputs one scalar prediction value for each time step, and compares it to the observation for that time step by computing a loss function, which in this case was the root-mean-square error (RMSE) between the observed and predicted discharges. As in Feng et al. (2020a), the Adadelta algorithm, an adaptive learning rate scheme (Zeiler, 2012), was selected as the optimization method for performing stochastic gradient descent on the model parameters of the neural network.

Normalization of inputs and outputs is a useful procedure to facilitate parameter updates by gradient descent. Normally, the loss function is defined over a mini-batch: the model is trained on many basins over the CONUS, and a random subset of hydrographs from some basins are put together to calculate the loss function. In this setup, however, wetter or larger basins contribute more to the loss function than the drier or smaller ones. To prevent this imbalance, we first normalized the daily streamflow by its area and mean annual precipitation to get a dimensionless streamflow, i.e., the runoff ratio, as the age-variable. Next, the distributions of daily streamflow and precipitation were transform at to e as close to a Gaussian distribution as possible, using the equation

$$v^* = \log_{10}(\sqrt{v} + 0.1) \tag{9}$$

where v is the original value and  $v^*$  is the transformed value. Finally, a standard transformation was applied to all the inputs by subtracting the CONUS-scale mean value and then dividing by the CONUS-scale standard deviation. The statistics used for normalization of the test period data were the same as those calculated for the training period data.

There were four hyperparameters: (i) the mini-batch size, which is the number of hydrographs that are put together a calculate the loss function before performing a weight update; (ii) the length of the hidrographs used for training; (iii) the number of hidden units, which is a direct representation of the learning capacity of the LSTM network; and (iv) the dropout probability, which is the probability that a weight is set to 0. As in Feng et al. (2020a), a mini-batch size of 100, an LSTM sequence length of 365, a hidden size of 256, and a dropout rate of 1.5 we e selected to run the model. The network training is stochastic in nature. Also similar in the previous setup, all networks in this paper were trained with n = 6 different random seeds. Streamflow predictions resulting from the different random seeds were combined into an ensemble-average prediction. All evaluation metrics were reported for the ensemble-average streamflow, except for the final model transferability experiment (For these experiments detailed in section 2.4.4, we could clearly reach the conclusion from one-random-seed experiments, so there was no need for multiple random seeds). All experiments were

implemented using adaptations from the PyTorch library (Paszke et al., 2017), and were performed on an NVIDIA GeForce GTX 1080 Ti GPU.

232

230

231

#### 2.2. Basin Datasets

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

233

Until now, there had not been a large-scale streamflow benchmark dataset containing extensive basins with reservoirs; CAMELS only has a small fraction of basins To compile such a dataset, we collected attributes, forcings, and stream low late for 3557 basins from GAGES-II, which also encompasses most of the CAME (see section 2.4). We selected 30 static physical attributes which fit into six egories: (1) basic identification and topographic characteristics, (2) percentages of an cover in the watershed, (3) soil characteristics, (4) geological characteristics, (5) local and cumulative dam variables, and (6) other disturbance variables (see Table S2 n Appendix for more details). Figure 1 plots ve ttributes of all basins including slope, forest the location of all 3557 sites and show fraction, soil permeability, normal stor of dams, and freshwater withdrawal. Basin mean forcing data for the period 01/01, 990–12/31/2009 was generated using the same method as for the CAMELS dataset, which is done by mapping a daily, gridded meteorological dataset, , 2016) to the chosen basin polygons. The Daymet dataset Daymet Version 3 (The me was acquired from he Cogle Earth Engine (GEE) data catalog (Gorelick et al., 2017) in the form of gridded es males of daily weather variables for the United States from 01/01/1980 to the present. The pasin mean daily time series forcing data were also obtained in GEE using auce functions. Pixels of the gridded data were determined to be in a region g to weighted reducers. Pixels were included if at least 0.5% of the pixel was in the region; their weight was the fraction of the pixel covered by the region. Daily average streamflow was the target variable, for which data for all gauges was downloaded from the USGS website (USGS, 2019). It should be noted that the Daymet data use UTC time (Spangler et al., 2019), while USGS daily values are based on local time (Sauer, 2002). It is difficult to correct this error as they were given in a daily format in the raw data. In this paper,

we directly use daily data from the Daymet dataset and the USGS to keep consistent with the CAMELS dataset, as many other studies did. Ideally, one would download sub-daily values from the USGS Instantaneous Values API and shift them to UTC before aggregating to days (or, vice versa, use an hourly forcing product and shift it to local time), as was done in some recent work (Gauch et al., 2020). While we do not think this error changes our conclusions, it calls attention to the need for revisions in datasets like CAMELS.

We also trained and tested models on the CAMELS dataset to allow for omparison to previous results. The CAMELS dataset (Addor et al., 2017; Newman et al., 2015, only included basins which experienced minimal human disturbance, noted as a convening gages, and excluded basins where human activities including artificial diversions reservoirs, and other activities in the basin or the channels significantly affected the natural flow of the watercourse (Falcone, 2011).



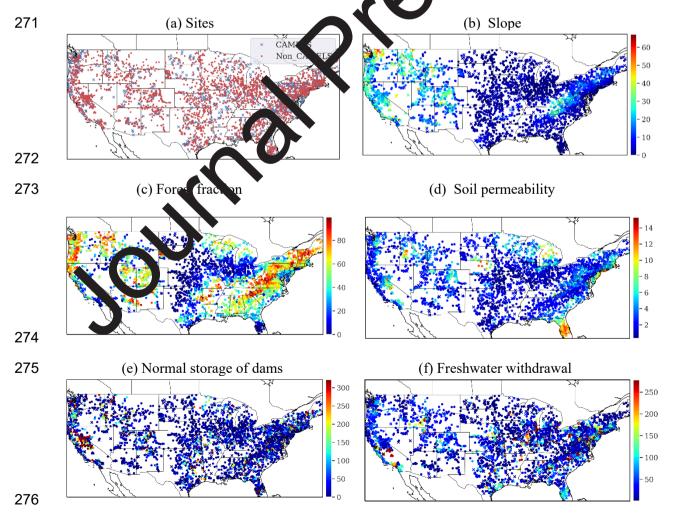


Figure 1. The location of all 3557 sites and characteristics of the corresponding basins. (a) Locations of all 3557 sites. Blue "x" markers are used to represent sites belonging to the CAMELS dataset, while red "o" points are the other, non-reference sites; (b) Slope: basin mean slope, as a percentage; (c) Forest fraction: percentage of basin with land cover "forest"; (d) Soil permeability: basin average permeability, inches/hour; (e) Normal storage of dams: total normal reservoir storage volume in a basin, megaliters of total storage per so km; (f) Freshwater withdrawal, megaliters per year per sq km. We excluded some scremely lirge values of (e) and (f) by choosing values below the 95% percentile value, in ride to more clearly show basin diversity.

# 2.3. Reservoir-related basin characteristics

Degree of regulation (*dor*) refers to the dimensive upstream reservoir storage as a percentage of the average streamflow, and is an important indicator of the impact of reservoirs on streamflow (Lehner et al., 2011). In the present study, it was calculated as the capacity-to-runoff ratio of a basin, defined as follows:

$$dor = \frac{nor}{\overline{q}} \tag{10}$$

where *nor* represents the sum of normal capacity of all reservoirs in a basin (m³ per km²), and  $\bar{q}$  is the estimated watershed mean annual runoff, or total volume of water annually leaving the basin the streamflow (m³ per km²), from GAGES-II. A *dor* value of 0.1 was set as the cutoff unit between basins with relatively little human regulation (small-*dor* basins) and basins with relatively large human regulation (large-*dor* basins) based on our preliminary analysis of the distribution of whole-CONUS model's performance across different basins as a function of *dor*. The *dor* is analogous to the commonly used metric of storage ratio (McMahon et al., 2007). A basin with *dor*=0.1 has the approximate storage of about a month of streamflow, which typically would be expected to have significant impact on daily streamflow yet is not enough to heavily modulate flow across seasons. On a side note, *dor* was not the threshold used by

CAMELS to select basins. CAMELS contains 344 small-dor basins and 32 large-dor basins, which represent a much smaller fraction of the CAMELS basins as compared to the overall CONUS.

We hypothesized that reservoir characteristics such as their purposes could be useful. To obtain these attributes, dams listed in the National Inventory of Dams (NID) database (US Army Corps of Engineers, 2018) were spatially joined with the boundary polygons of the basins. To minimize the influences of these differences on our results, we excluded a hich did not have matching dams included in NID and GAGES-II. Next, for ever the reservoir's normal capacity associated with each dam purpose culated. The purpose with the largest associated capacity was considered to be in major purpose of the collective dams in the basin. If there were more than one purpose sharing the largest capacity, we calculated normal storages of these purposes in order a importance (indicated by the order of the letters symbolizing the dam's purposes, ".ga"SC" indicates a primary purpose of the chose the most important purpose with the water supply followed by flood control) largest capacity. If still more than on propose was obtained, we treated them as being of equal importance, meaning that here were multiple main dam purposes listed for that basin. There were only a few basin with two categories of main dam purposes (only 1 basin had the Introl", and only 7 basins had the main dam purpose of main dam purpose of Debris "Navigation"), which was not enough to determine statistical characteristics, so they were excluded from the tatistical analysis. After all of these processing steps were complete, 656 e 5557-basin dataset were excluded from the statistical analysis in section 2.4.2: basins for a not have dams, 38 basins do not have dams listed in either the GAGES-II 610 or NID database, and 8 basins have main dam purposes of "Debris Control" or "Navigation". As a result, 2901 basins with 10 main dam purposes (Table 1) were available to analyze the influence of reservoir types (Table 2).

We added flags to describe the presence of water diversion, based on remarks and comments included in the GAGES-II dataset. "WR\_REPORT\_REMARKS" reported remarks pertinent to hydrologic modifications from the Annual Data Report (ADR) citation of the USGS,

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

and "SCREENING\_COMMENTS" reported screening comments from National Water-Quality Assessment (NAWQA) personnel regarding evidence of human alteration of flow, based on visual (primarily Google Earth) screening. We manually read through the text in these columns, and if there was some description with "diversion" or "divert" for a basin, the presence of diversion for this basin was regarded as "True"; otherwise it was assumed "False". Unfortunately, there was no available data regarding the volume of diversion, and hence diversion could only be used as a qualitative flag for our statistical analysis.

Table 1. Major reservoir purposes for basins in our dam characteristic ataset

Type	Purpose	Number of Basins
С	Flood Control and Stormwater Management	313
F	Fish and Wildlife Pond	94
Н	Hydroelectric	196
I	Irrigation	328
0	Other	163
Р	Fig Protection, Stock, or Small Farm Pond	66
R	Recreation	1207
S	Water Supply	426
	Tailings	52
X	Unknown	66

# 2.4. Experiments

#### 2.4.1. Temporal generalization tests

As we first wanted to determine the level of performance that could be achieved using one model over all 3557 basins in the full dataset (Table 2), an LSTM-based model (LSTM-CONUS) was trained and tested over all of these basins. For comparison to previous studies using the CAMELS dataset, we selected 523 basins (Table 2) from CAMELS (LSTM-CAMELS) to form a training set. The choice of 523 was made for multiple reasons. Firstly, the 3557basin dataset does not actually contain all of the CAMELS basins. In addition data from the GAGES-II dataset and the forcing data used in this study, Daynot Version 3 in GEE (last access in this study: 18 January 2020), were not exactly s those used for CAMELS. Finally, by removing some basins with large basin are there is a 531-basin subset of CAMELS which has often been selected as the ben himark set for rainfall-runoff modeling in previous work (Feng et al., 2020a; Kratzert et al., 2020a; Kratzer 19b). An intersection between the 3557 basins and this 531 benchmark CAME S deset basins resulted in the 523-basin All Nodels were trained using data from 1 January "baseline" CAMELS dataset we used he 1990 through 31 December 1999, and sting was done using data from 1 January 2000 through 31 December 2009.

361

362

363

364

365

366

367

368

369

370

371

372

360

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

# 2.4.2. Exploring the in parts reservoir attributes on model performance

There are tany reservoir attributes that could potentially inform improvements in streamflow models at such as dam storage or distance from gage location to dam. As the first paper (5 th, best of our knowledge) to study continental-scale streamflow prediction in damned mains in a deep learning context, we explored the impacts of multiple reservoir attributes and anthropogenic factors (details in Appendix Figure S2). Then, within the scope of this paper and partially consistent with McManamay (2014), we examined three major factors having significant influence on our model performance: capacity-to-runoff ratio (degree of regulation, *dor*), main dam purpose, and presence of diversion. As the models utilized in this study were basin-centric, these factors needed to be aggregated to each basin, which was done following the procedures discussed in Section 2.3.

Table 2. Datasets used in the this study

Name	Number of basins	Explanation
full dataset	3557	Basins with complete streamflow records during 1990/01/01-2009/12/31, selected from GAGES-II (section 2.4.1)
523-CAMELS dataset	523	Basins contained both in full datases and CAMELS (section 2.4.1)
dam characteristics dataset	2901	Subset of full dataset, containing basins used to explore the impacts of the three includes: capacity-to-runoff ratio ( <i>dor</i> ), dam purpose, and diversion (section 2.4.2)
zero- <i>dor</i> dataset	610	Subset of full data set, sor aining basins without dams (section 2, 3, 2, .4)
small- <i>dor</i> dataset	1762	Subset of full decases, containing basins with 0 < dor < 0. (section 2.4.3, 2.4.4)
large- <i>dor</i> dataset	1185	Subjet of full dataset, containing basins with <i>dor</i> ≥ 0.1 (section 2.4.3, 2.4.4)

# 2.4.3. Stratification by reservoh regime vs. pooling data together

379 F
380 k
381 p
382 k
383 a
384 c
385 t
386 v
387 f

For DL models in general providing more data often leads to model improvements. From the perspective of machine learning, then, lumping all data together would thus seem to be the obvious procedure to follow, given the likely beneficial impacts on modeling performance as we has simple implementation. However, it remains possible that stratification by rese voir attributes might result in clear separation basins with different latent (unknown) attributes. Hence, our research question 2 raised in the Introduction became two subquestions: (2A) Should we group all basins together, or classify basins into certain types and train models for each class separately to achieve the best performance? (2B) Do basins with varied reservoir regimes (no reservoir, small reservoir, or large reservoirs) function fundamentally differently? This could be proven true if basins trained in one regime cannot apply to basins in another regime.

To answer question 2A, all basins in the full dataset were divided into three groups (Table 2): zero-dor basins (dor=0), small-dor basins (0<dor<0.1) and large-dor basins (dor≥0.1). We trained models on these different groups individually, as well as together in various combinations. First, we trained and tested three LSTM-based models, called LSTM-Z, LSTM-S, and LSTM-L (we used "LSTM-x" to represent the LSTM-based models, which was different from the naming method for the datasets), on zero-dor, small-dor, ind-large-dor basins, respectively. Second, basins from two of the three groups were combined into training sets for three additional LSTM-based models: LSTM-ZS (trained on zero-dor and small-dor datasets), LSTM-ZL (trained on zero-dor and large-dor datasets), and large-dor datasets), but these three models were useful on basins from each of zero-dor, small-dor, and large-dor datasets. Finally, the testing esults of basins in these three groups were compared to results for the same balling from the LSTM-CONUS (trained on full dataset) model.

# 2.4.4. Model transferability experiments

To answer question (2B) vises in 2.4.3, we ran a set of predictions in ungauged basins (PUB) experiments, in valid medels trained in one set were tested in other sets. Further, when a model is trained as some basins and tested in others, the performance will naturally degrade. Therefore, we adoubt control experiments where models were trained and tested on the same categories or basins, which helped to disentangle the effects of reservoir regime and spatial extractal plants.

or example, zero-dor basins were divided into two batches (Train-z and PUB-z) with a ratio of 1:1 for training and test, respectively. We ensured that each of these cases was representative of the full group by including basins from every LEVEL-II ecoregion (Omernik and Griffith, 2014). The model trained on the Train-z set is then tested on Train-z itself, PUB-z and a subset (PUB-s) of the small-dor basins. These three test sets represent temporal generalization alone, spatial extrapolation and "spatial extrapolation+difference in reservoir

regime", respectively. Similarly, we separated the small-dor dataset into Train-s and PUB-s, and the large-dor dataset into Train-I and PUB-I. We also ran experiments with a mixed training set, e.g., Train-z and Train-s were merged to form one training dataset called Train-zs. Once trained on Train-zs, the LSTM-based model was tested individually on PUB-z and PUB-s. Two more training sets, combining zero-dor basins with large-dor ones (Train-zI), and pairing small-dor basins with large-dor ones were set up in the same way (Train-sI). It was not practical to attempt all possible combinations, but the combinations used sufficiently preswered the question (2B).

Finally, a fourth sub-experiment was added for comparison, a structuransferability of the LSTM-based model trained on the 523-CAMELS dataset. The Jasins of the 523-CAMELS dataset were also divided into the training (Train-c) modest (PUB-c). Then, the models trained on Train-c were tested on itself and other subsets (PUB-c/PUB-z /PUB-s/PUB-I). The details of all four of these sub-experiments are listed in Table 3.

Table 3. A summary of the training and testing datasets for sub-experiments exploring PUB with dams. All models were trained from January 1990 through December 1999, and tested from January 2000 through Pec abbe 2009. Multiple basin counts are given for each case of the first three sub-experiment, as we ran two tests (and therefore performed the basin groupings twice) for each case. For example, in the first sub-experiment, Train-z had 299 basins for the first run, and 309 basins for the second run. We list the Train-z and PUB-z datasets two each the first and second sub-experiments, because they belong to two independs a sub-experiments.

sub-experiment ID	training dataset (explanations)	test dataset (explanations)
1	Train-z (299/309 randomly selected	Train-z (same as the training set)
	zero- <i>dor</i> basins)	PUB-z (309/209 zero- <i>dor</i> basins that are different from those in Train-z)

		PUB-s (300/292 randomly selected small- <i>dor</i> basins)
	Train-zs (A mixture of 544/560 zero-dor or small-dor basins)	PUB-z (280/272 zero- <i>dor</i> basins that are different from those in Train-zs)
		PUB-s (280/272 small-dor basins that are different from those in Train-zs)
2	Train-z (295/305 randomly selected zero- <i>dor</i> basins)	Train-z (same as the falming set)
		PUB-z (3) 5/295 zero- <i>dor</i> basins that are different from those in Train-z).  NUB-I 297/289 randomly
		ela red large-dor basins)
	Train-zl (A mixture of 512/528 z to-vr or large-dor basins)	PUB-z (264/256 zero- <i>dor</i> basins that are different from those in Train-zl)
		PUB-I (264/256 large- <i>dor</i> basins that are different from those in Train-zl)
3	Train-s (871/ 79 ra Jomly selected	Train-s (same as the training set)
	s II-d(r.b.isins)	PUB-s (879/871 small- <i>dor</i> basins that are different from those in Train-s)
		PUB-I (639/634 randomly selected large- <i>dor</i> basins)
50	Train-sl (A mixture of 876/888 small-dor or large-dor basins)	PUB-s (444/438 small- <i>dor</i> basins that are different from those in Train-sl)
		PUB-I (444/438 large- <i>dor</i> basins that are different from those in Train-sl)
4	Train-c (257/264 basins in the 523- CAMELS dataset)	Train-c (same as the training set)
		PUB-c (264/257 basins that are different from the Train-c dataset, but still in the 523-CAMELS

#### dataset)

PUB-z (383 zero-dor basins that are different from the 523-CAMELS dataset)

PUB-s (1482 small-dor basins that are different from the 523-CAMELS dataset)

PUB-I (1169 large-dor busins that are different from the 52 CAMELS dataset)

439

440

441

442

443

444

445

446

447

448

449

450

451

452

#### 2.5. Metrics

In this study, the metrics used to mathematically quantify the racy of the models included bias, Pearson's correlation (Corr), the Nash-Sur liffe model efficiency coefficient (NSE) (Nash and Sutcliffe, 1970) and Kling-Gupta efficier (NSE) (Gupta et al., 2009). Bias is the mean difference between modeled and objected values. Corr is the linear correlation coefficient between modeled and observed values, and is not influenced by bias. NSE is a normalized statistic that determines the ative magnitude of the residual variance compared to the measured data variance. E is a nonlinear combination of correlation, flow variability measure, and bias; it is another a muon metric to evaluate how well the models perform. We the top 2% high flow volume range (FHV) and the percent also reported the percent be bias of the bottom 30% w flow volume range (FLV) (Yilmaz et al., 2008). FHV and FLV highlight the performance of the model for peak flows and baseflow, respectively. Metrics for in this study are reported for the test period (01/01/2000-12/31/2009). all experimen

453

454

455

#### 3. Results and Discussion

#### 3.1. CONUS-scale model with reservoirs

456

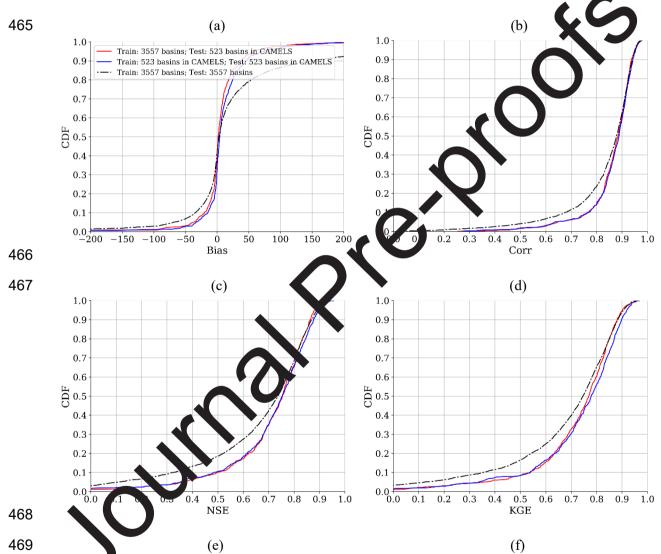
457

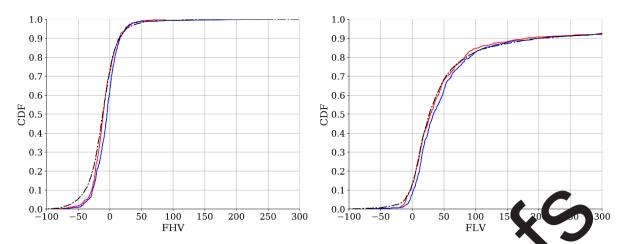
458

For the 3557 basins in the full dataset, the ensemble median NSE of the CONUS-scale model reached 0.74 (Figure 2c, details of ensemble experiments recorded in Appendix Table

S3). This value is at the same level as the previous benchmarks with the CAMELS reference-basin dataset (Feng et al., 2020a; Kratzert et al., 2020), despite that 83% of the 3557 basins have dams present in GAGES-II. When the models trained on CAMELS (LSTM-CAMELS) and CONUS (LSTM-CONUS) were tested on the 523-CAMELS baseline reference dataset, both achieved a median NSE values of 0.75 (Figure 2c, more details in Appendix Table S3).







470

471

472

473

474

475

Figure 2. Comparison of the empirical cumulative distribution functions CDr for the 523 basins tested in LSTM-CONUS and LSTM-CAMELS, and the 3557 fasin in LSTM-CONUS. The CDF of FLV does not reach 1.0 because the 30% low flot interaction for some basins is completely composed of zero-flow observations. Therefore for these basins, the percent bias is infinite, and thus the x-axis cannot include them.

476

477

478

479

480

481

482

483

484

485

486

487

488

489

The high NSE for the entire see newhat unexpected, because we had earlier thought that reservoirs would create allenges for LSTM and there may not be reliable mapping relationships that could be learned on a large scale. Comparing our results to those reported in the literature, a 1ST 10.74 certainly represents a state-of-the-art prediction for basins with reservoirs and a much more operationally-reliable model. Besides the values reported in literator summarized in the Introduction and Table S1, many of which reported negative NSTs for his challenging problem, the closest value we can find in the literature was (al. \ 008), who added reservoirs into a simple lumped hydrologic model, tested this Payan ( modern 40 pasins (mostly in France), and reported a mean NSE of 0.68. We would also like to note that the meteorological data for CONUS seems to have larger error than the European counterpart, which could lead to our model presenting an even higher NSE with European basins if we were to train our models there. In line with this hypothesis, some of our previous work showed that we could obtain a NSE of 0.84 for CAMELS-GB (Coxon et al., 2020), which

has 670 basins from United Kingdom (Ma et al., 2021), while the same model with the same training procedure could only achieve a NSE of 0.74 for CAMELS over CONUS.

492

493

494

495

496

497

498

499

500

501

502

503

504

505

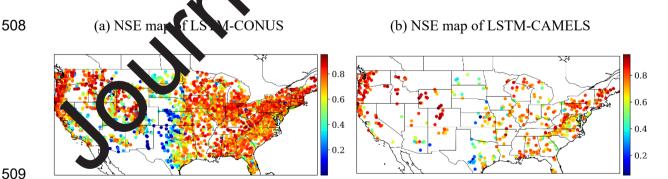
506

490

491

When tested on the 523-CAMELS dataset, the expanded dataset led to slightly improved overall bias with almost the same correlation but slightly decreased KGE (noticeable by comparing red and blue lines in Figure 2a-b.d). Since KGE is a composite metric of correlation, flow variability, and bias, we suspect that additional samples in the enlarged the flow variability, which makes it a little more difficult for LSTM the flow variability for the 523 basins. This hypothesis can be further by looking at the values for FHV and FLV. The median FHV values when tested on a 3 CAMELS basins were -10% for LSTM-CONUS and -4% for LSTM-CAMELS, showing a minor increase in highflow bias for the expanded dataset (Figure 2e). In contrast for the expanded dataset (Figure 2e). In contrast for the expanded dataset (Figure 2e). simulations were improved by the use of a bigger rain dataset, as the median FLV values N.1-CAMELS (Figure 2f). Compared to CAMELS, were 28% for LSTM-CONUS, and 33% a higher fraction of basins with large reservoirs we suspect the expanded set may co which attenuate the peak flow, and hence the LSTM-CONUS model tended to predict lower peaks.

507



509

510

511

Figure 3. NSE spatial patterns of the ensemble results of (a) LSTM-CONUS and (b) LSTM-CAMELS.

LSTM-CONUS and LSTM-CAMELS both showed good performance in the northwestern CONUS and most parts of the eastern CONUS, but had relatively poor performance on the Great Plains, Texas, Oklahoma, Kansas, and parts of California (Figure 3). The regional distribution of NSEs is largely in line with earlier work (Feng et al., 2020a), where basins on the Great Plains and the extremely-dry southwestern border performed poorly with LSTM-based modeling. Evidently these basins in the central CONUS continue to pose challenges for LSTM despite the larger dataset, perhaps because the larger still large basins where the homogeneous assumption of the LSTM-based models bleak down.

521

513

514

515

516

517

518

519

520

### 3.2. Analysis of the impacts of reservoir-related factors

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

522

Using the results from the CONUS-scale simulation (L M-CONUS), we explored the uncertainty of the current LSTM-based model guided by three attributes: the capacity-to-runoff e of the dam and its associated reservoir, and the ratio (degree of regulation, dor), the pur as a clear pattern regarding dor: regardless of the presence of diversion (Figure 4a). The expurpose, the overall model performance, as quantified by the median NSE, was always better for small-dor basins than for larger-or ones (see Figure 4d). This observation differs from red with a process-based model (Shin et al., 2019), which previously-reported results ob had more difficult predicting the streamflow of basins with small-capacity reservoirs (corresponding to man dor). The management policies of reservoirs could change over time nink that is potentially the reason why the model did not perform as well for large-dor and we ever, for small-dor reservoirs, the model still delivered excellent performance so basin Langes in policies may not have resulted in dramatic impacts for these small reservoirs. A first-order visualization of the impacts of other control variables are given in Appendix Figure S2.

538

539

540

Exploring model uncertainty based on dam purpose not only showcased the uncertainty of the LSTM-based models, but also clearly indicated that different types of

reservoirs exert varied influences on streamflow. Among all the various dam purposes, basins with reservoirs mainly for recreation (R) or water supply (S) were easier to model. It may be inferred that the water storages of these reservoirs changed relatively little on a daily scale to achieve their purposes and therefore had less impact on the streamflow than other reservoirs (Ryan et al., 2020). Three types of reservoir purposes stood out as being more challenging to predict (Figure 4b); fire protection or farm ponds (P), irrigation (I), and hydroelectric (HL Basins with "P" reservoirs, for any dor value range regardless of the presence of difficult to predict and had the worst performance of all those in the small and ate, ory. This indicates that LSTM had trouble finding a universal relationship to more es for a chain of many small, individually-regulated ponds. Difficulty in modeling frig tion reservoirs was not unexpected, as it has been shown that irrigation water usage has specific seasonal variations, and is related to the crop type, field, and other site-specific formation (Shin et al., 2019). Critical information that would help with modeling for these basins, such as water use and the operational policies of hydroelectric (H) dams timing, is not generically available. Like seek to optimize electricity production, appliare therefore influenced by the prices on the local electricity grid (Giuliani et al., 2014), which were not included in this dataset.

substantially decreased NSE values (Figure 4a). For The presence of divers η at there are smaller NSE values for dam purposes "I", "O", instance, it is visibly ar par nt "P", and "R" in the basis with diversion. This was also expected: diversion influences the water balance, but because no information about the quantity of diverted water was available TIM pased model, the model couldn't understand the imbalance, leading to reduced to the L formance. A clearer separation is seen in the results of four specific cases, which prediction. combinations of only two categorical variables -- the dor value range, and the presence of diversion (Figure 4c). The median NSEs for small-dor basins without diversion, small-dor basins with diversion, large-dor basins without diversion and large-dor basins with diversion were 0.78, 0.76, 0.65, and 0.62, respectively. It was evident that LSTM could reach the best performance in small-dor basins without diversion, while the worst performance

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

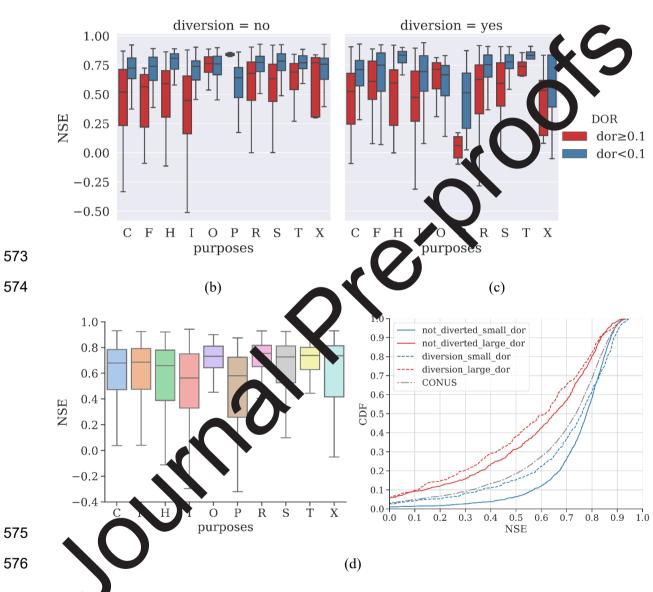
564

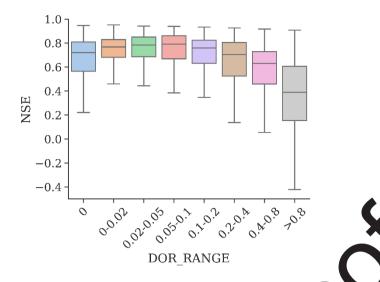
565

566

occurred in large-*dor* basins with diversion, and thus the effects of the two factors seem to be additive.

572 (a)





577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

Figure 4. (a) NSE distributions with three categorical variables: dol value range ("small-dor"

basins have 0<dor< 0.1 and "large-dor" basins have dor ≥ purposes of reservoirs in a basin, and presence of diversion. Dam purposes are C: A od Control and Stormwater Management; F: Fish and Wildlife Pond: H. Hyan ectric: I: Irrigation: O: Other: P: Fire Protection, Stock, or Small Farm Policy, R. Recreation; S: Water Supply; T: Tailings; and X: Unknown. (b) NSE distribution ar basins with different main dam purposes. (c) NSE empirical cumulative distribution function curses from LSTM-CONUS and four cases resulting from riables: dor range and presence of diversion. The blue and combinations of two categorical va represent the NSE distributions of small-dor basins with and without green lines respect ely diversion, which vere picked out from the ensemble result of LSTM-CONUS. The red and tively indicate the NSE distributions of large-dor basins with and without orange grey dashed line represents the empirical CDF of LSTM-CONUS. (d) NSE as n of dor values all 3557 basins; the ranges of dor values: 0, (0, 0.02], (0.02, 0.05], (0.05, 0.1], (0.1, 0.2], (0.2, 0.4], (0.4, 0.8], >0.8, where "(]" means a left side half open interval;

the correspond numbers of basins in each range: 610, 1076, 377, 309, 311, 277, 247, 350;

other plots in this figure are for dam characteristics dataset shown in table 2.

594

The main challenges for LSTM-based modeling of reservoirs are clearly delineated (Figure 4a): LSTM had difficulty predicting streamflow for large-dor basins with dams for fish and wildlife, flood control, hydroelectric power generation, irrigation, and fire protection, with difficulty increasing in this order. Diversion further added to the challenge. To our knowledge, such identification of specific challenges has not been previously reported. Additionally, it was not previously clear that these challenges mainly exist only for large-dor basins. Small-dor basins, even those with reservoirs for irrigation and hydroelectric purposes, call the reasonably captured by LSTM, presumably because they have limited adaptive capacity. LSTM can approximate an optimal information extractor, which suggests that we are not apply sufficient information needed to model the more challenging cases and provides a trageted direction for future work.

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

595

596

597

598

599

600

601

602

603

604

605

dor is apparently a major control on LSAN model performance (Figure 4d). ro or basins, have the highest performance. The Interestingly, small-dor basins, instead median NSE in the 0.05-0.1 dor bin is alrost 0.8, a very high number (we offer explanations later). Below dor<0.1 human decisions cannot shift water availability across seasons. As discussed earlier, basins with do 10. Thave the reservoir storage equivalent to approximately w. As dor gets bigger than this amount, they have more one month of average su an capability to regulate flow on a seasonable scale, and the impact of human choice becomes more prominent. We also found the basin with more reservoirs could have equivalent or higher perform ince (Figure S2I), which suggests the difficulty may have mainly come from one or ams. Due to sometimes unpredictable human decisions influenced and also the marity in such decisions, e.g., shift in reservoir management policies, the *dor*>0.1 becomes increasingly difficult to simulate. This figure is also the basis for us to choose dor=0.1 as the threshold. Despite the challenges for large-dor basins, we nonetheless note that even for these basins, LSTM obtained a median NSE of 0.65 for basins without diversion, which is higher than many literature values reported in Table S1. To put things even further into context, a recent study for a basin with a major dam (USGS 11462500, Russian River near Hopland,

California, *dor* = 0.17) reported oftentimes negative daily NSE values and correlation between 0.5 to 0.8 for different months of the year (Kim et al., 2020). In contrast, the CONUS-scale model developed in this study reported a very high NSE value of 0.88 and correlation of 0.94 for this specific station. For a different comparison, the National Water Model reported an NSE of 0.62 for reference basins in CAMELS (Kratzert et al., 2019a).

628

629

627

623

624

625

626

#### 3.3. Impacts of training dataset

630

631

632

633

634

635

636

637

638

639

640

641

642

643

Our experimental results suggest that datasets with different anges can be trained together to enhance overall performance, and at the very last, bed training should not exert a significant detrimental impact on the model (Fig re 5 see more details in Tables S3 and S4, Appendix). With the inclusion of small-dor basins the training set (LSTM-ZS), there was a small improvement in predictions for unanimed basins (Wilcoxon signed-rank test:  $p=4.9 \times 10^{-6}$ ). For small-dor basing rere no clear differences in test performance when training with zero-dor basins together. In the large-dor basins, as compared to the result of LSTM-L (training with only large- or basins), all other cases reported slightly increased NSE values and fewer "catastrophic latures" (cases with NSE close to or smaller than 0), suggesting that new it ormat on was brought in by pooling information together. It is possible that the inclusion of zero-dor or small-dor basins allowed the model to better understand natural flows and el abled better modeling of the large-dor basins. Such a pattern fits with our eral obsertations obtained from training DL models.

644

645

646

647

648

649

We did see a slight exception to this pattern, however, when adding large-*dor* basins to the training set. When large-*dor* basins were added to the training set, a minute deterioration in NSE was observed when this model was tested on zero-*dor* and small-*dor* basins: the median NSE decreased from 0.72 to 0.71 for LSTM-ZL (left panel of Figure 5a, Wilcoxon signed-rank test:  $p=1.3 \times 10^{-4}$ ), and there was a declination from 0.79 to 0.78 shown for LSTM-

SL (center panel of Figure 5a, Wilcoxon signed-rank test:  $p=1.2 \times 10^{-32}$ ). We hypothesize that operations of large reservoirs are characteristically different from those of smaller reservoirs, and therefore the inclusion of large reservoirs introduced some noise to the data and made it more difficult for LSTM to grasp a universal pattern. Nevertheless, the adverse impact was quite minor. This result, along with our other observations of LSTM-CONUS (Section 3.1), also imply that it should be possible to fine-tune the LSTM-CONUS model for a local region to obtain refined simulations.

We were surprised to see that small-dor basins had notably ligher NSE values (median NSE ~0.79) than zero-dor basins (median NSE ~0.72) (Figure 5a. Two hypotheses could potentially explain this phenomenon: first, that the small-dor asins may be concentrated in certain areas, e.g., mountainous areas, where NSEs tend to be ligher; second, that a small-dor reservoir may serve as a buffer to boost the sterage of the system, thereby reducing the impacts of flash precipitation peaks which are or illerging to model (Feng et al., 2020a). Looking at the basins on a map are in the parameter space (Figure 5b), however, while mountainous basins do have higher NS s, the zero-dor and small-dor basins are mixed in space and there is no spatial aggregation of one or the other. Therefore, we reject the first hypothesis (concentration) and lead toward the second one (buffer).

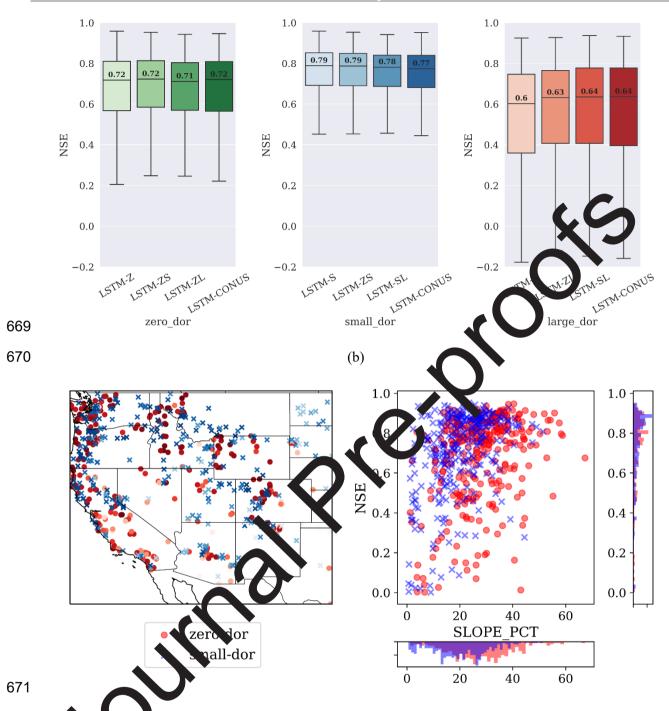


Fig. e . (a) Boxplots of NSE values for zero-dor basins (Z, dor=0), small-dor basins (S,

0 < do < 0.1) and large-dor basins (L, dor $\geq 0.1$ ). Green, blue, and red boxes show the results from models respectively tested on zero-dor, small-dor, and large-dor basins, while the training sets are noted on the x-axis labels. For each color, the lightest-colored box was trained solely with the same subset of basins on which it was tested, while the others had additional subsets included in the training sets. Basins in the test sets were always subsets of the training sets,

and the models were trained in 1990-1999 and tested in 2000-2009. (b) The left part is a NSE map of the western CONUS where small-dor and zero-dor basins coexisted. There are 303 zero-dor basins and 310 small-dor basins shown here. The right is a scatter plot of the relationship between NSE and SLOPE PCT (mean watershed slope, as a percent). The NSE values are part of the results for LSTM-CONUS (section 3.1). Red circular markers represent the zero-dor sites, and blue x-shaped markers represent the small-dor sites. For the map only. sites with lighter colors have lower NSE values.

685

678

679

680

681

682

683

684

686

687

688

689

690

691

Additionally, we were also surprised to see that LSTM sk onably good performance on even large-dor basins, with median NSE values .64 in the overall CONUS training sets (the rightmost boxplot in Figure 41), it spectively, which were still comparable to SAC-SMA's median NSE of 0.65 (Feng et a) 2 (va) for reference basins. This localing reservoirs as compared to earlier result suggests a large advantage of LSTM for methods.

692

# 3.4. The PUB experiments and model transferability

694

695

696

697

698

699

700

701

702

703

704

705

693

In the introduction, were the NSE values for dammed basins As we asked in the tio similar to previous esults with CAMELS because these basins in fact behave similarly? If this was not the case, New different are these basins? Our stratified PUB experiments showed that e substantial differences between zero-dor, small-dor, and large-dor basins such that there wa els trained only on one type of basin to other basin types caused significant apply perference drop that could not be explained solely by spatial extrapolation (Figure 6). For example, the median NSE values for "Train-z", "PUB-z", and "PUB-s" were 0.65, 0.51, and -0.06, respectively (Figure 6a). The scenario Train-z was a temporal test only, so this NSE value of 0.65 represents model performance without spatial extrapolation (this value was lower than LSTM-Z shown in Figure 5a because the training sample size was smaller: the zero-dor basins were randomly split for this experiment, as explained in section 2.4.4). The decline from

0.65 to 0.51 for PUB-z was then due to spatial extrapolation in the same zero-*dor* group. The more dramatic decline from 0.51 for PUB-z to -0.06 for PUB-s can be entirely attributed to the behavioral difference between zero-*dor* and small-*dor* basins. We also note larger declination for large-*dor* basins (Figure 6b-c), with median NSE values of -0.19 and 0.18 for the PUB-I cases.

Including diverse basins in the training dataset substantially elevated overall PUB performance. The mixed training sets (Train-zs, Train-zl, and Train-sl, the boxes on the light side of each panel in Figure 6a-c) had greatly improved median NSE values, as yell as greatly reduced incidences of catastrophic failures (cases with NSE close to an

It is noteworthy to mention that when we trained a model foler, or pasins subset from the 523-CAMELS dataset and then tested it on the other trasing of 523-CAMELS as well as zero-, small-, and large-*dor* basins, the model gave outrightelist strous results for PUB-z, PUB-s, and PUB-I (Figure 6d). This means that CAM £LS basins, as they are reference basins, differ fundamentally from the others, example from the zero-*dor* basins. This result distinctively highlights the danger of using CAMELS basins as the whole training set for continental-scale modeling, and also suggests we cannot simply ignore small reservoirs or simply treat them as being equivalent to reference basins.

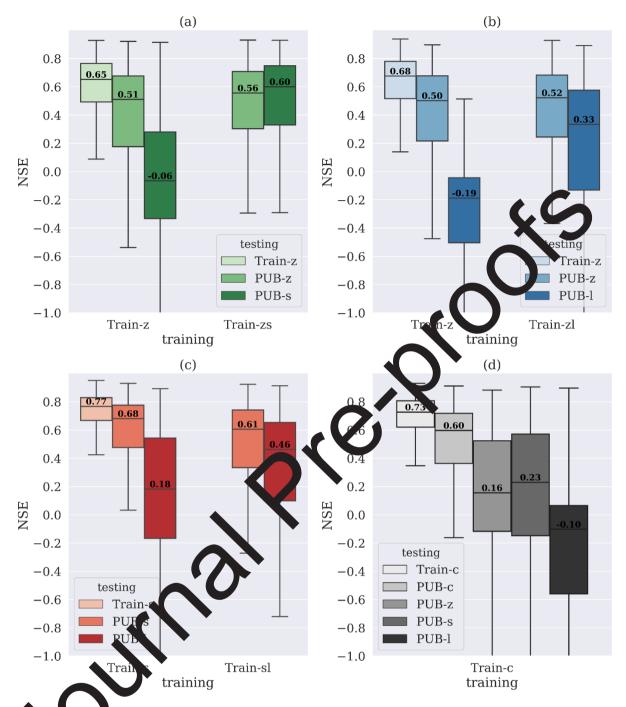


Figure & Boy plots from PUB sub-experiments where training and testing basins were from different combinations of basin types: c indicates 523-CAMELS, z indicates zero-dor basins, s indicates small-dor basins, and I indicates large-dor basins. Combinations of letters indicate that a combination of the indicated basin types were used (refer to Table 3 for details). The drop in performance from training basin-located test results to PUB-basin-located test results of the same type (e.g. Train-z vs PUB-z) represents the effect of spatial extrapolation, while the drop across different basin type combinations (e.g. PUB-z vs PUB-s) represents the effect

of migrating models across reservoir regimes. A side note: the PUB-c in (d), with a median of 0.60, is not comparable to other PUB tests in the literature. Here we only used ~260 CAMELS basins as training data and did not employ an ensemble for different random seeds (so as to be inline with other experiments in this figure). This test is solely shown to highlight the difference between the CAMELS basins and the others.

737

732

733

734

735

736

#### 3.5. Further Discussion

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

738

sed on input In future work, we could allow LSTM to estimate model una attributes, as shown in the modeling of soil moisture (Fang et I., ) and rainfall-runoff (Klotz et al., 2020). To further improve modeling capabilities for he more challenging cases, it could be useful to incorporate more information regarding was suse, electricity price patterns, and estimated diversion rates from sources like water management models (Yates et al., 2005) Gio ani et al., 2016). Fine-tuning may be another into the context of optimization process approach to improve predictions in Nor challenging basins (Sampson et al., 2020). For example, Ma et al. (2021) transferred their model trained on the CAMELS basins over to a few basins in Sichuan province To ha and obtained better results than the model trained with nated information such as distribution of the storage capacity all local basins. Other reservoil among the basin's eservirs, surface water area, or storage change in a basin may also be used as inputs though an encoder unit (Feng et al., 2020b). Moreover, physics-guided Read et al., 2019) could be employed to provide more stability where machin tata is scarce. In addition, a distributed version of the deep learning models could monk ring. It the spatial heterogeneity of a basin and may perform better than the lumped ones for large basins. In the future, machine-learning-based routing schemes (Bindas et al., 2020) can be added to support flood modeling in major rivers.

As a rule of thumb for DL models, pooling data together almost always helped improve modeling, which was confirmed by the zero-dor and small-dor cases shown in this study. However, here the large-dor basins could slightly pull down the metrics for other cases, which

deviated, albeit in a minor way, from this rule. We think that this was due to a combination of the rainfall-runoff processes from different basins having very dissimilar patterns, and the information from the inputs not being enough to discern differences between reservoir regimes, causing the LSTM-based model to struggle in fitting all of this information into one universal model. We suspect that the large-*dor* basins represent an extreme case of the problem of unmodelable dissimilarity in geoscience. The cut-off *dor* of 0.1 in this paper is an operational threshold, but may not be the only choice. Other *dor* cut-off values may also be applied ble, but this was not the focus of this paper. Future work should concentrate on low o incorporate more information and tune the model structure to train a universal result. for all non-regulated/regulated basins.

770

760

761

762

763

764

765

766

767

768

769

#### 4. Conclusion

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

771

s of modeling rainfall-runoff processes with Prior work has documented the LSTM in reference basins with minimal anthropogenic impacts. However, to our knowledge, no previous deep-learning baset study focused on basins significantly impacted by reservoir operations at a continental scale or the modeling implications of reservoir attributes. For this consisting of 3557 basins over the CONUS, and trained an work, we created a new LSTM-based mode which achieved an ensemble test median Nash Sutcliffe model efficiency coefficient (NSE) 10.74. This performance was at the same record level as reported for LS M-based modeling benchmarks, which showed for the first time that many previou be modeled as part of the standard basin rainfall-runoff and storage processes. mese results provide the first benchmarks for basins with and without reservoirs: zerodor, small-dor, and large-dor basin subsets had median NSE values of 0.72, 0.79, and 0.60, respectively. Furthermore, the NSE value for even the most challenging large-dor basins in the model over the CONUS (0.64) was still comparable to that of the current operational hydrologic model, SAC-SMA, trained and tested only with reference basins (0.65) (Feng et al.,

2020a), which further highlights the effectiveness of LSTM as a competitive option for emulating basins with reservoirs for large-scale hydrologic modeling.

Our results provided us with a coherent modeling strategy and some useful lessons. We showed that zero-dor and small-dor basins behave characteristically differently (and are also different from CAMELS reference basins), which strongly suggests that we cannot simply ignore smaller reservoirs out of convenience and treat them as natural flow, the standard practice in some process-based models. If using a data-driven model, the Jene cial strategy we determined for small reservoirs was to include their reservoir attributes and train a lumped, uniform model that simulated them as part of the basin off processes. We showed that basins with different dor values can be trained to get ver a large dataset to obtain record-level modeling performance, a strategy which could greatly simplify the modeling process. If using a process-based model, the process-based model, the process approach may be to modify parameters in the model, e.g., linear reser oir parameters, to represent the impacts of obtained the best performance in small-dor basins smaller reservoirs. The LSTM-based ma without diversion, especially for those vitt reservoirs for water supply and recreation. For the large-dor reservoirs of certain types, i.e., lire protection or farm ponds, hydroelectric, and irrigation dams which are most of culto model, we may adopt a mixed approach to represent M is already very strong with respect to feature extraction, them separately. Considering it is likely that more felevent information, e.g., electricity prices or irrigation water demand, will be needed to imply we their simulation. This paper is the first time such a systematic analysis previded from a data-driven perspective. has beg

damed basins: there must be sufficient representation of small-dor and large-dor basins in the training set. Dammed and undammed basins behave characteristically differently, and migrating models between them can be dangerous: when a model trained only on CAMELS reference basins or zero-dor basins was tested on basins with dams present, we encountered catastrophic failures. We showed that pooling all data together for model training tended to improve results, and even when it did not (likely due to insufficient input information and very

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

heterogeneous	training	data	bringing	in	noise),	the	inclusion	of	training	data	from	othe
scenarios still d	id not sig	ınifica	ntly jeopa	ardi	ize the r	esult	S.					

817

815

816

#### **Acknowledgements**

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

818

The lead author, Wenyu Ouyang, was supported by the China Scholarship Council for one year of study at the Pennsylvania State University. Chaopeng Shen w supported by National Science Foundation OAC #1940190. We thank the reviewers whose comments have helped to substantially improve the basins along with their dor and NSE values from the LSTM-CONUS N are provided as an attachment in the Appendix. Forcing data used in this st dy a e available from the Daymet dataset website (https://doi.org/10.3334/ORNLDAAS/1 (a); The GAGES-II dataset can be downloaded from the U.S. Geological Sur ev JSGS) GAGES-II website flow data can be downloaded from USGS Water (https://doi.org/10.3133/70046617); Str g/10.5066/F7P55KJN); Reservoir attribution data Data for the Nation website (http://dx. can be downloaded at National A ventory of Dams website (https://nid.sec.usace.army.mil) he CAMELS dataset can be downloaded at CAMELS from U.S. Army Corps of En ine D6G73C3Q) from the U.S. National Center for website (http://dx.doi.or Atmospheric Research. The entire project code is available at GitHub (https://github.com. QuyangWenyu/HydroSPDB). We appreciate the LSTM code at GitHub mhpi/hydroDL). Many thanks to Google Earth Engine regarding the data (https:// the forcing data used in this study. proce

837

838

### **Appendix**

839

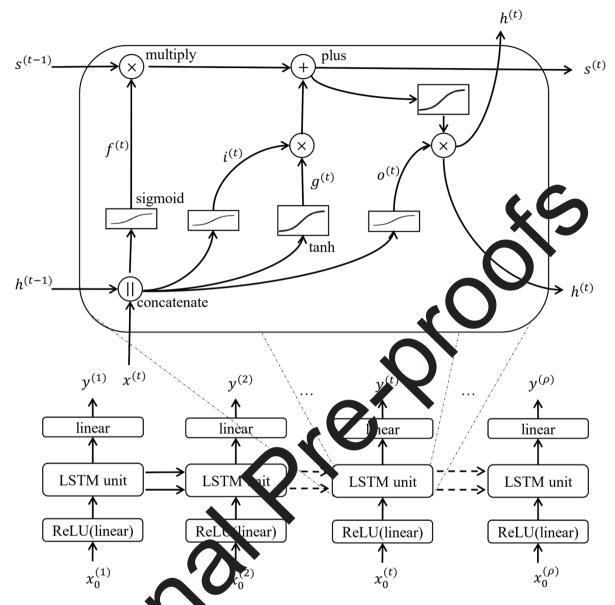
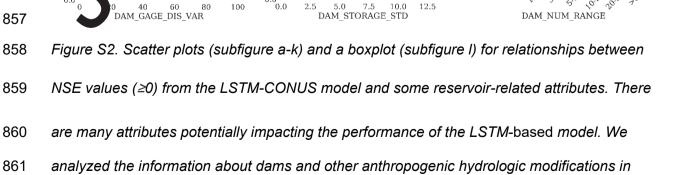


Figure S1. The illustration of the LSTM-based model structure and its unit.  $x_0^{(t)}$  is the vector of raw inputs for the line step t,  $\rho$  is the length of time sequence of LSTM in the training period. ReLU(litear, is the rectified linear unit,  $x^{(t)}$  is the input vector to the LSTM cell,  $g^{(t)}$  is the output of the impulsation,  $f^{(t)}$ ,  $i^{(t)}$ ,  $o^{(t)}$  are the forget, input and output gates, respectively,  $s^{(t)}$  represents the states of memory cells,  $h^{(t)}$  represents hidden states, and  $y^{(t)}$  is the predicted output which is compared to streamflow observations.

850 (a) NDAMS (b) STOR\_NOR (c) RAW\_DIS\_NEAREST\_MAJ\_DAM

#### 1.0 1.0 1.0 0.8 0.8 0.6 0.6 0.4 0.4 0.2 0.2 0.2 0.0 750 1000 1250 1500 NDAMS 1000 2000 3000 4000 5000 6000 STOR\_NOR 50 100 150 RAW\_DIS\_NEAREST\_MAJ\_DAM (e) FRESHW WITHDRAWAL (f) PCT IRRIG (d) RAW AVG DIS ALLDAMS 1.0 1.0 0.8 0.8 0.8 0.6 0.4 0.2 0.2 0.0 0.0 20 30 40 PCT\_IRRIG\_AG 0 100 150 200 250 RAW\_AVG\_DIS\_ALLDAMS 00 1000 1500 2000 FRESHW\_WITHDRAWAL 2500 (g) POWER SUM MW (h) PDEN BLOCK (i) ROADS KM SQ KM 1.0 1.0 0.8 0.8 0.0 NSE 0.6 0.4 0.2 0.2 0.2 0.0 0.0



0.2

2000 3000 PDEN\_BLOCK

(k) DAM STORAGE STD

4000

1.0

0.8 0.6

0.4 0.2 5.0 7.5 10.0 12.5 15.0 ROADS\_KM\_SQ\_KM

(l) DAM NUM RANGE

851

852

853 854

855

856

1.0

0.6

10000 15000 20000 25000 POWER\_SUM\_MW

(j) DAM GAGE DIS VA

862	the basin in the GAGES-II dataset. (a) NDAMS: number of dams in a basin; (b) STOR_NOR:
863	dam normal storage in a basin, megaliters total storage per sq km; (c)
864	RAW_DIS_NEAREST_MAJ_DAM: raw straight line distance of gage location to nearest
865	major dam in watershed, km. Major dams are defined as being >= 50 feet in height (15m) or
866	having storage >= 5,000 acre feet in GAGES-II; (d) RAW_AVG_DIS_ALLDAMS: raw
867	average straight line distance of gage location to all dams in watershed, km; (e)
868	FRESHW_WITHDRAWAL: freshwater withdrawal, megaliters (1000 cubic means) per year
869	per sq km; (f) PCT_IRRIG_AG: percent of watershed in irrigated agriculture; (g),
870	POWER_SUM_MW: sum of MW operating capability of electric generating plants in
871	watershed of type "coal", "gas", "nuclear", "petro", or "water"; (h) PDE LEZOCK: population
872	density in the watershed, persons per sq km; (i) ROADS_KM_S\_kM: road density, km of
873	roads per watershed sq km; (j) DAM_GAGE_DIS_VAR: the condicient of variation of the
874	distances from each dam to the gage location in a baum (k) DAM_STORAGE_STD: the
875	standard deviation (std) of the normal stances (stor) of reservoirs in a basin; we set
876	std(log(stor+1)) as the x-axis variable, to means the natural logarithm; (I)
877	DAM_NUM_RANGE: the ranges of dam numbers 0, 1, (1, 3], (3, 5], (5, 10], (10, 20], (20,
878	50], >50, where "(]" means a left lide half open interval; the correspond numbers of basins in
879	each range: 610, 362, 26, 28, 375, 442, 437, 619.
880	
881	
882	
883	Table \$1.25 servoir simulation results in the literature that do not use recent observations (i.e.
884	data milation or data integration). For comparison, our median NSE values reported here
885	were 0.74 for the whole set and 0.78 for basins with small reservoirs. For comparability, we
886	did not include papers that used continual inputs of recent observations of inflow, outflow, or
887	storage

Reference	Metric	Description
Shin et al.	No NSE reported. Correlations of monthly	high-resolution continental-scale reservoir
(2019)	outflow for the new scheme	scheme (grid-centric) which improved the
	$(R_{new})$ ranged from -0.07 to 0.63 with a median of 0.25,	simulations of reservoirs greatly over the
	which were higher than Hanasaki et al. (2006) and Biemans et al. (2011)	contiguous United States. Tested over six
	schemes.	reservoirs in the Missouri, caramento,
		Columbia, San Joaquin, and Colorado River
		Basins
Voison et al.	Best monthly NSE of	An improved grid-centric re ervoir formulation to
(2013)	regulated flow is 0.62. Negative NSEs for two	the heavily domined Columbia River Basin.
	other locations	Authors shawer performance metrics for
		moran, regulated flow at three locations.
Wu and Chen	NSE of outflow ≈ 0.3	reservoir operation scheme to decide outflow
(2012)		and its distribution on hydropower, water supply
		and impoundment purposes according to the
		inflow and storage. Authors calibrated the
	~~~	coefficients of equations in the new scheme
		during 1965-1984 and validated the scheme in
•		the period 1987-1988 for the Xinfengjiang
	<b>J</b>	reservoir
VO		
Kim et al., (2020)	Positive monthly NSEs of daily runoff discharges for	A grid-centric scheme inside the NWM. Tested
	real scheduled release;	on four locations and 21 hydrographs. An NSE
	most simulated releases brought negative NSEs	of 0.78 was reported for a short period (~11
	(reading off Figure 7)	days) of hourly simulation at one of the locations,
		but Figure 7 showed mostly negative NSEs.

Zajac et al. Best NSE of streamflow is 0.61 (reading off Figure 6) (2017)

Global daily streamflow simulations of a spatially distributed LISFLOOD hydrological model in 390 stations during 1980-2013

Zhao et al. NSE of 0.74 and 0.51 for (2016) outflow of two reservoirs, respectively.

multi-purpose reservoir module with predefined complex operational rules was integrated into the Distribute Soil Vegetation Model (DHSVM) rformance of the model was tested ov upper Brazos River Basin in Texas, where two reservoirs, Lake Whitney and Lake, are located

Payan et al.

Mean NSE = 0.68

(2008)

46 basins nos v in France).

The dailty of the meteorologic dataset in the US, used in this dataset, is potentially lower than the European counterpart. Our work showed that we could obtain NSE=0.84 for CAMELS-GB (Coxon et al., 2020), which has 670 basins from United Kingdom (Ma et al., 2021), while the same model with the same training procedure can achieve only 0.74 for CAMELS over CONUS, consistent with other studies. Beck et al., (2020) also showed that NSE for US basins are not higher than global basins.

Dang et al.

(2020)

A NSE range of 0.68-0.79 for the calibration period, but no value was reported for the validation period

A novel variant of VIC's routing model to simulate the storage dynamics of water reservoirs for the Upper Mekong. However, this study focused on

the effect of parameter compensation during calibrating the model or without the reservoir module. Hence, the author did not report on the test period.

889

890

891

892

Table S2. Summary of the forcing and attribute variables used as inputs the STM-based

893 model

Variable T	ype	Variable Name	Descriptio	Unit
Forcing		dayl	Day length	s
		prcp	Precipitation	mm/day
		srad	Solar radiation	W/m2
			Snow water equivalent	mm
			Maximum temperature	°C
		tmin	Minimum temperature	°C
	W.	vp	Vapor pressure	Pa
Attroutes	Basic identification	DRAIN_SQKM	Watershed drainage area	km <sup>2</sup>
J	and topographic characteristic	ELEV_MEAN_M_ BASIN	Mean watershed elevation	m
	S	SLOPE_PCT	Mean watershed slope	%
		STREAMS_KM_S Q_KM	Stream density	km of streams per watershe d km <sup>2</sup>

	Percentages of land cover in the watershed	DEVNLCD06	Watershed percent "developed" (urban)	-
		FORESTNLCD06	Watershed percent "forest"	-
		PLANTNLCD06	Watershed percent "planted/cultivated" (agriculture)	-
		WATERNLCD06	Watershed percent Open Water	5
		SNOWICENLCD0 6	Watershed percent Perennial Ice/Snow	-
		BARRENNLCD06	Watershed percels Natural Barren	-
		SHRUBNLCD06	Watershild percent Shrubland	-
		GRASSNLCD06	Valurshed percent Horba teous (grassland)	-
		WOODYVE NL D06	Watershed percent Woody Wetlands	-
		EMERGWL TNLC	Watershed percent Emergent Herbaceous Wetlands	-
	Soil characteristic s	AV CAVE	Average value for the range of available water capacity for the soil layer or horizon	inches of water per inches of soil depth
•		PERMAVE	Average permeability	inches/h
10		BDAVE	Average value of bulk density	g/cm <sup>3</sup>
2		ROCKDEPAVE	Average value of total soil thickness examined	inches
	Geological characteristic s	GEOL_REEDBU SH_DOM	Dominant (highest percent of area) geology	-
		GEOL_REEDBU SH_DOM_PCT	Percentage of the watershed covered by the dominant geology type	-

Cl	ocal and umulative am variables	NDAMS_2009	Number of dams in watershed	-
		STOR_NOR_200 9	Dam storage in watershed ("NORMAL_STORAGE")	megaliter s/km²
		RAW_DIS_NEAR EST_MAJ_DAM	Raw straightline distance of gage location to nearest major dam in watershed.	km
di	other isturbance ariables	CANALS_PCT	Percent of stream kilometers coded as "Canal", "Ditch", "Pipeline"	5
		RAW_DIS_NEAR EST_CANAL	Raw straightline distance of gage location to near st canal/ditch/pipeline in watershed	km
		FRESHW_WITH DRAWAL	Freshwa er wit drawal megaliters på year per sqkm	1000 m <sup>3</sup>
		POWER_SUM_M W	Cam of operating capability of electric generation power plants in watershed of type "coal", "gas", "nuclear", "petro", or "water"	MW
		PDL \_2000_BLO	Population density in the watershed	persons/ km²
•	"	O/DS_KM_SQ KM	Road density	km of roads per watershe d km <sup>2</sup>
. 0	J	IMPNLCD06	Watershed percent impervious surfaces	%

Table S3. Detailed ensemble results of LSTM-based models in this study

Model	Section in the "Experiments"	Random seed	NSE median	Ensemble NSE median
	Exportitionio			modian

		Journal Pre-p	proofs	
LSTM-CONUS	2.4.1	123	0.71	0.74
		1234	0.71	
		12345	0.72	
		111	0.69	
		1111	0.72	
		11111	0.71	
LSTM-CAMELS	2.4.1	123	0.73	0.6
		1234	0.74	
		12345	0.74	
		111	0.74	
		1111	0.68	
		11111	0.18	
LSTM-Z	2.4.3	123	J.69	0.72
		1234	0.65	
		95 5	0.71	
		111	0.69	
		1111	0.70	
	1	11111	0.68	
LSTM-S	2.4.3	123	0.77	0.79
10,		1234	0.77	
2		12345	0.78	
•		111	0.78	
		1111	0.76	
		11111	0.76	

		Journal Pre- <sub>r</sub>	proofs	
LSTM-L	2.4.3	123	0.52	0.60
		1234	0.58	
		12345	0.57	
		111	0.54	
		1111	0.59	_ <b>k</b> S
		11111	0.59	
LSTM-ZS	2.4.3	123	0.76	
		1234	0.74	
		12345	0 5	
		111	0.76	
		1111	0.77	
		1/1	0.76	
LSTM-ZL	2.4.3	123	0.64	0.66
		1234	0.63	
	.(1)	12345	0.64	
	<b>)</b> ,	111	0.63	
10		1111	0.64	
3		11111	0.63	
LSTM-SL	2.4.3	123	0.72	0.75
		1234	0.73	
		12345	0.72	

111	0.72
1111	0.72
11111	0.72

Table S4. Ensemble testing results of basins with different dor ranges in affect models

904 (Section 3.3)

1			
sub-experiment ID	Test basins (number of basins)	Training models	median NSE
1	zero- <i>dor</i> basins (610)	LSTHY	0.72
		LSI 2S	0.72
		LSTM-ZL	0.71
		LSTM-CONUS	0.72
2	small- or vasins	LSTM-S	0.79
. (		LSTM-ZS	0.79
		LSTM-SL	0.78
10		LSTM-CONUS	0.77
3	large- <i>dor</i> basins (1185)	LSTM-L	0.60
		LSTM-ZL	0.63
		LSTM-SL	0.64
		LSTM-CONUS	0.64
·			

908

#### References

909

913

914

915

916

917 918

919

920

921

922

923

924

925

926

927 928

929 930

931 932 933

934

935

936

941

942 943 944

945

946 947

948

949

950

951

952

953 954

955

956

- 910 Addor, N., Newman, A.J., Mizukami, N., Clark, M.P., 2017. The CAMELS data set: 911 catchment attributes and meteorology for large-sample studies. Hydrol. Earth Syst. 912 Sci. 21, 5293-5313. https://doi.org/10.5194/hess-21-5293-2017
  - Ayzel, G., Kurochkina, L., Kazakov, E., Zhuravlev, S., 2020. Streamflow prediction in ungauged basins: benchmarking the efficiency of deep learning, in: E3S Webof Conferences. EDP Sciences, p. 01001. https://doi.org/10.1051/e3sconf/202016301001
  - Beck, H.E., Pan, M., Lin, P., Seibert, J., van Dijk, A.I.J.M., Wood, E.F., 2020. Glo Distributed Parameter Regionalization Based on Observed Stream Jow Headwater Catchments. J. Geophys. Res. Atmospheres 125, e20 9JD0 1485. https://doi.org/10.1029/2019JD031485
  - Biemans, H., Haddeland, I., Kabat, P., Ludwig, F., Hutjes, R.W.A. Heinke, J., von Bloh, W., Gerten, D., 2011. Impact of reservoirs on river discharge, and reservoir supply during the 20th century. Water Resour. Res. 47. https://doi.org/10.1029/2009WR008929
  - ough the river network utilizing Bindas, T., Shen, C., Bian, Y., 2020. Routing flood waves to physics-guided machine learning and the Muskingum-Lunge Method, in: American Geophysical Union (AGU). Presented at the AGU Fall Meeting 2020, American inge Method, in: American Geophysical Union (AGU).
  - Coxon, G., Addor, N., Bloomfield, J.P., Freer, F., Hannaford, J., Howden, N.J.K., Lane, R., Lewis, M., Robinson, Wyener, T., Woods, R., 2020. CAMELS-GB: hydrometeorological time series and landscape attributes for 671 catchments in Great Britain. Earth Syst. Sci. Lat. 12, 2459–2483. https://doi.org/10.5194/essd-12-2459-2020
  - Dang, T.D., Chowdhury, A.F.M. Galelli, J., 2020. On the representation of water reservoir storage and operations in la re-scale hydrological models: implications on model parameterization and clip, te change impact assessments. Hydrol. Earth Syst. Sci.
  - 24, 397–416. https://do.org/10.5194/hess-24-397-2020
    Ehsani, N., Fekete, B.M., Vörö Carty, C.J., Tessler, Z.D., 2016. A neural network based general reserver operation scheme. Stoch. Environ. Res. Risk Assess. 30, 1151– 1166. https://doi.org/10.1007/s00477-015-1147-9
  - SAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow
  - (Report). A. stori, VA. https://doi.org/10.3133/70046617
    Fang, K., Kife, D., Jawson, K., Shen, C., 2020. Evaluating the Potential and Challenges of n D. Sehainty Quantification Method for Long Short-Term Memory Models for Soil Moistu e Predictions. Water Resour. Res. 56, e2020WR028095. //doi.org/10.1029/2020WR028095
  - Shen, C., 2020. Near-Real-Time Forecast of Satellite-Based Soil Moisture Using ong Short-Term Memory with an Adaptive Data Integration Kernel. J. Hydrometeorol. 21, 399-413. https://doi.org/10.1175/JHM-D-19-0169.1
  - Feng, D., Fang, K., Shen, C., 2020a. Enhancing Streamflow Forecast and Extracting Insights Using Long-Short Term Memory Networks With Data Integration at Continental Scales. Water Resour. Res. 56, e2019WR026793. https://doi.org/10.1029/2019WR026793
  - Feng, D., Lawson, K., Shen, C., 2020b. Prediction in ungauged regions with sparse flow duration curves and input-selection ensemble modeling. ArXiv Prepr. ArXiv201113380.
  - Gal, Y., Ghahramani, Z., 2016. Dropout as a Bayesian approximation: representing model uncertainty in deep learning, in: Proceedings of the 33rd International Conference on

- 959 International Conference on Machine Learning - Volume 48, ICML'16. JMLR.org, 960 New York, NY, USA, pp. 1050-1059.
  - Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., Hochreiter, S., 2020. Rainfall-Runoff Prediction at Multiple Timescales with a Single Long Short-Term Memory Network. Hvdrol. Earth Svst. Sci. Discuss. 2020, 1–25. https://doi.org/10.5194/hess-2020-540
    - Giuliani, M., Castelletti, A., Pianosi, F., Mason, E., Reed, P.M., 2016. Curses, tradeoffs, and scalable management: Advancing evolutionary multiobjective direct policy search to improve water reservoir operations. J. Water Resour. Plan. Manag. 142, 04015050. https://doi.org/10.1061/(ASCE)WR.1943-5452.0000570
    - Giuliani, M., Herman, J.D., Castelletti, A., Reed, P., 2014. Many-objective reservoir policy identification and refinement to reduce policy inertia and myopia in water management, Water Resour, Res. 50, 3355-3377. https://doi.org/10.1002/2013WR014700
    - Gochis, D., Barlage, M., Dugger, A., FitzGerald, K., Karsten, L., McAllister, M. Mills, J., RafieeiNasab, A., Read, L., others, 2018. The WRF-Hydra technical description, (Version 5.0). NCAR Tech. Note 107.
    - Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., Mo 7. Google Earth Engine: Planetary-scale geospatial analysis for everyone. Bit Remote. Sensed Data Tools Appl. Exp. 202, 18–27. https://doi.org/10.1016 (.rst. 2017.06.031)
    - Grill, G., Lehner, B., Thieme, M., Geenen, B., Tickner, D., Antone, F., Babu, S., Borrelli, P., Cheng, L., Crochetiere, H., Ehalt Macedo, H., Filgur, Fas, R., Soichot, M., Higgins, J., Hogan, Z., Lip, B., McClain, M.E., Meng, J., Mulligan, M., Iilsson, C., Olden, J.D., Opperman, J.J., Petry, P., Reidy Liermann, C., Sároz, , Salinas-Rodríguez, S., Schelle, P., Schmitt, R.J.P., Snider, J., Tan T. Tockner, X., Valdujo, P.H., van Soesbergen, A., Zarfl, C., 2019. Mapping ne world's free-flowing rivers. Nature 569, 215–221. https://doi.org/10.1038/s4158-0-11-1-9
    - Gupta, H.V., Kling, H., Yilmaz, K.K., Maxin Z, F., 2009. Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. J. Hydrol. 377, 80–31. Jups://doi.org/10.1016/j.jhydrol.2009.08.003

      Gutenson, J.L., Tavakoly, A.A., Wahl, M.D., Follum, M.L., 2020. Comparison of generalized non-data-driven lake and reservoir fouting models for global-scale hydrologic
    - forecasting of reservoir outling at diurnal time steps. Hydrol. Earth Syst. Sci. 24, 2711–2729. https://doi.org/10.3194/hess-24-2711-2020

      Hanasaki, N., Kanae, S., Okr, 7., 2506. A reservoir operation scheme for global river routing models. J. Hydrol. 27, 2506. https://doi.org/10.1016/j.jhydrol.2005.11.011

    - Hochreiter, S., 1998, T. e Val shing Gradient Problem During Learning Recurrent Neural Nets and Poblet Solutions. Int. J. Uncertain. Fuzziness Knowl.-Based Syst. 06, 107-116 https://doi.org/10.1142/S0218488598000094
    - iter, S., Schmidwuber, J., 1997. Long Short-Term Memory. Neural Comput. 9, 1735–1780. Hps: Idoi.org/10/bxd65w Hochreiter, S
    - International Rivers, 2007. Damming Statistics [WWW Document]. Int. Rivers. URL ttps://archive.internationalrivers.org/damming-statistics
    - a, L., Johnson, L.E., Gochis, D., Cifelli, R., Han, H., 2020. An experiment on Kim, servoir representation schemes to improve hydrologic prediction: coupling the ational water model with the HEC-ResSim. Hydrol. Sci. J. 65, 1652–1666. https://doi.org/10.1080/02626667.2020.1757677
    - Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Klambauer, G., Hochreiter, S., Nearing, G., 2020. Uncertainty Estimation with Deep Learning for Rainfall-Runoff Modelling. ArXiv E-Prints arXiv:2012.14295.
- 1008 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A.K., Hochreiter, S., Nearing, G.S., 2019a. 1009 Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine 1010 Learning, Water Resour, Res. 55, 11344–11354. 1011 https://doi.org/10.1029/2019WR026065
- 1012 Kratzert, F., Klotz, D., Hochreiter, S., Nearing, G.S., 2020. A note on leveraging synergy in 1013 multiple meteorological datasets with deep learning for rainfall-runoff modeling.

963 964

965

966

967

968

969

970

971

972

973

974

975

976

977

982 983

984

994 995

996

997

998

999

1000 1001

1002

1003

1004

1005

1006

- 1014 Hydrol. Earth Syst. Sci. 2020, 1–26. https://doi.org/10.5194/hess-2020-221
- Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., Nearing, G., 2019b.
   Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets. Hydrol. Earth Syst. Sci. 23, 5089–5110.

1018 https://doi.org/10.5194/hess-23-5089-2019

- Lawrence, D.M., Fisher, R.A., Koven, C.D., Oleson, K.W., Swenson, S.C., Bonan, G.,
  Collier, N., Ghimire, B., van Kampenhout, L., Kennedy, D., others, 2019. The
  Community Land Model version 5: Description of new features, benchmarking, and
  impact of forcing uncertainty. J. Adv. Model. Earth Syst. 11, 4245–4287.
  https://doi.org/10.1029/2018MS001583
- Lehner, B., Liermann, C.R., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P., Döll, P., Endejan, M., Frenken, K., Magome, J., Nilsson, C., Robertson, J.C., Rodel R. Sindorf, N., Wisser, D., 2011. High-resolution mapping of the world's recryolis and dams for sustainable river-flow management. Front. Ecol. Environ. 9, 494, 502. https://doi.org/10.1890/100125
  - Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., Sharma, L., Shin, C., 2021. Transferring hydrologic data across continents -- leveraging data rich, gions to improve hydrologic prediction in data-sparse regions. Water Resource. Res. n/a, e2020WR028600. https://doi.org/10.1029/2020WR028600
  - McMahon, T.A., Pegram, G.G.S., Vogel, R.M., Peel, M.C., 2007. Levisiting reservoir storage–yield relationships using a global streamflor / dat bace. Adv. Water Resour. 30, 1858–1872. https://doi.org/10.1016/j.advwatres. 007. 2.003
  - McManamay, R.A., 2014. Quantifying and generalizing hydron acresponses to dam regulation using a statistical modeling approach. J. Hydro. 519, 1278–1296. https://doi.org/10.1016/j.jhydrol.2014.08\_0 3
  - Mulligan, M., van Soesbergen, A., Sáenz, L., 2020, GOODD, a global dataset of more than 38,000 georeferenced dams. So: Lata 1–8. https://doi.org/10.1038/s41597-020-0362-5
  - Nash, J.E., Sutcliffe, J.V., 1970. River for forecasting through conceptual models part I A discussion of principles, J. Hydro, 10, 282–290. https://doi.org/10.1016/0022-1694(70)90255-6
  - Newman, A.J., Clark, M.P., Sampsen, K., Wood, A., Hay, L.E., Bock, A., Viger, R.J., Blodgett, D., Brekke Z., John, J.R., Hopson, T., Duan, Q., 2015. Development of a large-sample watershed-scare hydrometeorological data set for the contiguous USA: data set characters ics and assessment of regional variability in hydrologic model performance it vdrol. Farth Syst. Sci. 19, 209–223. https://doi.org/10.5194/hess-19-209-2015
  - Omernik, J.M., Criffic, G.E., 2014. Ecoregions of the conterminous United States: evolution of a bierarc ical spatial framework. Environ. Manage. 54, 1249–1266. https://loi.org/10.1007/s00267-014-0364-1

  - Patterson, L.A., Doyle, M.W., 2018. A Nationwide Analysis of U.S. Army Corps of Engineers Resour. Assoc. 54, 543–564. https://doi.org/10.1111/1752-1688.12622
  - Payan, J.-L., Perrin, C., Andréassian, V., Michel, C., 2008. How can man-made water reservoirs be accounted for in a lumped rainfall-runoff model? Water Resour. Res. 44. https://doi.org/10.1029/2007WR005971
- Read, J.S., Jia, X., Willard, J., Appling, A.P., Zwart, J.A., Oliver, S.K., Karpatne, A., Hansen, G.J., Hanson, P.C., Watkins, W., others, 2019. Process-guided deep learning predictions of lake water temperature. Water Resour. Res. 55, 9173–9190. https://doi.org/10.1029/2019WR024922
- Ryan, J.C., Smith, L.C., Cooley, S.W., Pitcher, L.H., Pavelsky, T.M., 2020. Global Characterization of Inland Water Reservoirs Using ICESat-2 Altimetry and Climate

- 1069 Reanalysis. Geophys. Res. Lett. 47, e2020GL088543. 1070 https://doi.org/10.1029/2020GL088543
- 1071 Sampson, A.K., Hale, E., Lambl, D., 2020. Big Data for Specific Places in Hydrologic 1072 Modeling, in: American Geophysical Union (AGU), Presented at the AGU Fall 1073 Meeting 2020. American Geophysical Union (AGU).
  - Sauer, V.B., 2002. Standards for the Analysis and Processing of Surface-Water Data and Information Using Electronic Methods (Report No. 2001–4044), Water-Resources Investigations Report. https://doi.org/10.3133/wri20014044
  - Shen, C., 2018. A transdisciplinary review of deep learning research and its relevance for water resources scientists. Water Resour. Res. 54, 8558-8593. https://doi.org/10.1029/2018WR022643
  - Shin, S., Pokhrel, Y., Miguez-Macho, G., 2019. High-Resolution Modeling of Release and Storage Dynamics at the Continental Scale. Water Reso 787-810. https://doi.org/10.1029/2018WR023025
  - Spangler, K.R., Weinberger, K.R., Wellenius, G.A., 2019. Suitability of grid Lea datasets for use in environmental epidemiology. J. Expo. Sci. Env. on. E idemiol. 29, 777-789. https://doi.org/10.1038/s41370-018-0105-2
  - Thornton, P.E., Thornton, M.M., Mayer, B.W., Wei, Y., Devarakonda R., V se, R.S., Cook, R.B., 2016. Daymet: Daily Surface Weather Data on a 14 n r North America, Version 3. ORNL Distributed Active Archive Center. https://doi.org/10.3334/ORNLDAAC/1328
  - Turner, S.W.D., Doering, K., Voisin, N., 2020. Data-Driven oir Simulation in a Largeeser Scale Hydrological and Water Resource Model. Water Resource Res. 56.
  - e2020WR027902. https://doi.org/10.1029/2020Wx027902 US Army Corps of Engineers, 2018. National inventory of dams [WWW Document]. URL https://nid.sec.usace.army.mil/
  - ete. Web interface [WWW Document]. U. S. USGS, 2019. National water information
  - Geol. Surv. URL https://water\_ata.u\_gs.gc\_v/nwis?

    Voisin, N., Li, H., Ward, D., Huang, N., Wamosta, M., Leung, L.R., 2013. On an improved sub-regional water resources management representation for integration into earth system models. Hydrol. Earth Syst. Sci. 17, 3605-3622. https://doi.org/10.5194/hess-17-3605-2013
  - Wu, Y., Chen, J., 2012. An Operation Based Scheme for a Multiyear and Multipurpose Reservoir to Enhance Mac oscale Hydrologic Models. J. Hydrometeorol. 13, 270–283. https://doi.org/10.105/JHM-D-10-05028.1
  - r, I., 2020. A Rainfall-Runoff Model With LSTM-Based Sequence-to-Xiang, Z., Yan, J., Del Sequence arm. v. Water Resour. Res. 56, e2019WR025326. https://dei.org/10.1029/2019WR025326
- Yang, S., Yang, D. Chen, J., Zhao, B., 2019. Real-time reservoir operation using recurrent neural cetw rks and inflow forecast from a distributed hydrological model. J. Hydrol. 1107 1108 1109 24229. https://doi.org/10.1016/j.jhydrol.2019.124229
- .. Razavi, S., Elshamy, M., Davison, B., Sapriza-Azuri, G., Wheater, H., 2019. 1110 sentation and improved parameterization of reservoir operation in hydrological 1111 1112 nd land-surface models. Hydrol. Earth Syst. Sci. 23, 3735–3764. 1113 https://doi.org/10.5194/hess-23-3735-2019
- 1114 Yates, D., Sieber, J., Purkey, D., Huber-Lee, A., 2005. WEAP21—A demand-, priority-, and 1115 preference-driven water planning model: part 1: model characteristics. Water Int. 30, 1116 487-500. https://doi.org/10.1080/02508060508691893
- Yilmaz, K.K., Gupta, H.V., Wagener, T., 2008. A process-based diagnostic approach to 1117 1118 model evaluation: Application to the NWS distributed hydrologic model. Water 1119 Resour. Res. 44. https://doi.org/10.1029/2007WR006716
- 1120 Zajac, Z., Revilla-Romero, B., Salamon, P., Burek, P., Hirpa, F.A., Beck, H., 2017. The 1121 impact of lake and reservoir parameterization on global streamflow simulation. J. 1122 Hydrol. 548, 552-568. https://doi.org/10.1016/j.jhydrol.2017.03.022
- 1123 Zeiler, M.D., 2012. ADADELTA: An Adaptive Learning Rate Method. CoRR abs/1212.5701.

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085 1086

1087

1088

1089

1090

1091

1092 1093 1094

1095

1100

1101 1102 1103

1104 1105

1124 1125 1126 1127	Zhao, G., Gao, H., Naz, B.S., Kao, SC., Voisin, N., 2016. Integrating a reservoir regulation scheme into a spatially distributed hydrological model. Adv. Water Resour. 98, 16–31. https://doi.org/10.1016/j.advwatres.2016.10.014
1128	WO ran the experiments, produced the visualizations, and wrote the initial manuscript, DF
1129	provided assistance in modeling, KL, DF, LY, CZ and CS edited the manuscript. CS
1130	conceived the study and revised the manuscript.
1131	4.6
1132	
1133	Highlights
1134 1135 1136 1137 1138 1139 1140	<ol> <li>LSTM achieved state-of-the-art performance for modeling basins with reservoirs.</li> <li>Reservoir types, capacity-to-runoff ratio (<i>dor</i>) and diversion contra-streamflow.</li> <li>LSTM performed well for basins with reservoirs that store about a month of flow.</li> <li>It is crucial to include basins with reservoirs in the thining set.</li> <li>Large-<i>dor</i> basins with certain types of dams are more difficult for LSTM.</li> </ol>