# Regionalization in a global hydrologic deep learning model: from physical descriptors to random vectors

Xiang Li<sup>1,1,1</sup>, Ankush Khandelwal<sup>1,1,1</sup>, Xiaowei Jia<sup>2,2,2</sup>, Kelly Cutler<sup>1,1,1</sup>, Rahul Ghosh<sup>1,1,1</sup>, Arvind Renganathan<sup>1,1,1</sup>, Shaoming Xu<sup>1,1,1</sup>, J L Nieber<sup>1,1,1</sup>, Christopher J Duffy<sup>3,3,3</sup>, Michael Steinbach<sup>1,1,1</sup>, Vipin Kumar<sup>4,4,4</sup>, and Kshitij Tayal<sup>5,5</sup>

<sup>1</sup>University of Minnesota <sup>2</sup>University of Pittsburgh <sup>3</sup>Pennsylvania State University <sup>4</sup>Department of Computer Science and Engineering, University of Minnesota <sup>5</sup>University of Minnesota Twin Cities

November 30, 2022

#### Abstract

Streamflow prediction is a long-standing hydrologic problem. Development of models for streamflow prediction often requires incorporation of catchment physical descriptors to characterize the associated complex hydrological processes. Across different scales of catchments, these physical descriptors also allow models to extrapolate hydrologic information from one catchment to others, a process referred to as "regionalization". Recently, in gauged basin scenarios, deep learning models have been shown to achieve state of the art regionalization performance by building a global hydrologic model. These models predict streamflow given catchment physical descriptors and weather forcing data. However, these physical descriptors are by their nature uncertain, sometimes incomplete, or even unavailable in certain cases, which limits the applicability of this approach. In this paper, we show that by assigning a vector of random values as a surrogate for catchment physical descriptors, we can achieve robust regionalization performance under a gauged prediction scenario. Our results show that the deep learning model using our proposed random vector approach achieves a predictive performance comparable to that of the model using actual physical descriptors. The random vector approach yields robust performance under different data sparsity scenarios and deep learning model selections. Furthermore, based on the use of random vectors, high-dimensional characterization identifies the uniqueness of catchments, thereby improving regionalization performance in gauged basin scenario when physical descriptors are uncertain, or insufficient.

# Regionalization in a global hydrologic deep learning model: from physical descriptors to random vectors

## Xiang Li<sup>1</sup>, Ankush Khandelwal<sup>2</sup>, Xiaowei Jia<sup>3</sup>, Kelly Cutler<sup>2</sup>, Rahul Ghosh<sup>2</sup>, Arvind Renganathan<sup>2</sup>, Shaoming Xu<sup>2</sup>, Kshitij Tayal<sup>2</sup>, John Nieber<sup>1</sup>, Christopher Duffy<sup>4</sup>, Michael Steinbach<sup>2</sup>, Vipin Kumar<sup>2</sup>

<sup>1</sup>Department of Bioproducts and Biosystems Engineering, University of Minnesota Twin Cities, St.Paul, MN, USA

<sup>2</sup>Department of Computer Science and Engineering, University of Minnesota Twin Cities, Minneapolis, MN, USA

<sup>3</sup>School of Computing and Information, University of Pittsburgh, Pittsburgh, PA, USA

<sup>4</sup>Department of Civil and Environmental Engineering, Pennsylvania State University, State College, PA,

USA

## 13 Abstract

1

2

3

4

6

7

8

9

10

11

12

Streamflow prediction is a long-standing hydrologic problem. Development of mod-14 els for streamflow prediction often requires incorporation of catchment physical descrip-15 tors to characterize the associated complex hydrological processes. Across different scales 16 of catchments, these physical descriptors also allow models to extrapolate hydrologic in-17 formation from one catchment to others, a process referred to as "regionalization". Re-18 cently, in gauged basin scenarios, deep learning models have been shown to achieve state 19 of the art regionalization performance by building a global hydrologic model. These mod-20 els predict streamflow given catchment physical descriptors and weather forcing data. 21 However, these physical descriptors are by their nature uncertain, sometimes incomplete, 22 or even unavailable in certain cases, which limits the applicability of this approach. In 23 this paper, we show that by assigning a vector of random values as a surrogate for catch-24 ment physical descriptors, we can achieve robust regionalization performance under a 25 gauged prediction scenario. Our results show that the deep learning model using our pro-26 posed random vector approach achieves a predictive performance comparable to that of 27 the model using actual physical descriptors. The random vector approach yields robust 28 performance under different data sparsity scenarios and deep learning model selections. 29 Furthermore, based on the use of random vectors, high-dimensional characterization im-30 proves regionalization performance in gauged basin scenario when physical descriptors 31 are uncertain, or insufficient. 32

#### **1 Introduction**

In hydrology, streamflow prediction is essential for the forecast of water supply, floods, 34 and droughts. It is a challenging task because of interacting hydrological processes (Beven, 35 1989, 1987; Freeze & Harlan, 1969; Freeze, 1974), spatial-varying parameter uncertain-36 ties (Beven & Binley, 1992), and limited observations (Blöschl & Sivapalan, 1995). These 37 challenges have motivated the advancement of hydrologic models from simple to com-38 plex. Encompassing more underlying hydrological processes, a complex hydrologic model 39 includes more hydrologic parameters and detailed catchment physical descriptors to ad-40 dress the complexities (Beven, 2001, 2002) and associated scaling issues (McDonnell et 41 al., 2007). But parameterizing such a complex hydrologic model for any individual catch-42 ment becomes difficult when hydrologic data are unavailable. Thus, regionalization, which 43 is defined as "how to extrapolate hydrologic information from one area to another" (Blöschl 44 & Sivapalan, 1995), specifies a research topic of modeling catchment runoff prediction 45 using hydrologic information from multiple catchments, which will be given a brief back-46 ground review in section 2. 47

Regionalization heavily relies on physical descriptors, such as, soil porosity, catch ment elevation, etc. These physical descriptors account for hydrologic complexities and
 regional differences and are thus intensively used in regionalized hydrologic models, ei ther process-based or data-driven.

Recently, Kratzert et al. (2019a) have presented a regionalized data-driven hydro-52 logic model that greatly outperforms local process models. Specifically, they trained a 53 single deep learning model (LSTM, abbreviated for the Long Short-Term Memory net-54 works) for 531 basins in the US CAMELS (Catchment Attributes and Meteorology for 55 Large Sample studies) dataset (Addor et al., 2017) and show that it is able to greatly 56 outperform the well-established process-based models (e.g., SAC-SMA (Burnash, 1995), 57 VIC (Liang et al., 1994), etc) that have been individually parameterized for each basin, 58 and thus offer a better route to regionalization (Kratzert et al., 2019a). 59

Building such a model requires streamflow observation and weather forcings for many 60 basins with diverse physical descriptors. It also relies upon the fact that all relevant basin 61 physical descriptors are available and of high quality. Performance of such models may 62 suffer if some of the descriptors are missing or are incorrect/uncertain. Our paper presents 63 an approach where it is possible to build a data driven regionalized model even in the 64 absence of any basin specific physical descriptors. It is able to use the weather forcing 65 and streamflow data from a set of basins to build a global model without having any in-66 formation about the physical descriptors of individual basins (For the background in-67 formation of the global model, please see section 2). However the structure of this model 68 is identical to the one used by Kratzert et al. (2019a), as it only replaces the individual 69 catchment physical descriptors by random vectors that simply provide a distinct char-70 acterization to each basin. Our results show that this approach provides global models 71 at least as good as the ones produced using the knowledge of all available physical de-72 scriptors. But the performance is much better relative to the scenario where some of phys-73 ical descriptors are missing and/or are incorrect/uncertain. 74

We note that the random vector and physical descriptor approaches are not in confict and in fact give comparable results. In fact, for ungauged basins, Kratzert et al.'s model can be used (Kratzert et al., 2019b) and shows that physical descriptors serve as a bridge between gauged basins and ungauged basins. In our approach, the random vectors do not connect gauged basins and ungauged basins due to the lack of streamflow observation for the ungauged basins.

The paper is organized as follows. Section 2 introduces relevant background infor-81 mation, in particular the regionalization. Section 3 explains the details of the random 82 vector method as well as the deep learning architecture involved. This section also ex-83 plains the dataset and the set up of the experiment. The experiment includes an exhaus-84 tive analysis on the applicability of our proposed random vector methods under various 85 data scarce situations and modeling structures. Section 4 lists our benchmarking results 86 and the exhaustive analysis of the random vector applicability. Section 5 highlights sci-87 entific implications from our results and suggests a few future directions. Section 6 sum-88 marizes the scientific conclusions. 89

## 90 2 Background

Performing hydrologic prediction from multiple catchments, regionalization is closely
related to the problem addressed in "prediction in ungauged basins" (PUB) (Sivapalan
et al., 2003), and most literature uses "PUB" and "regionalization" interchangeably (Pagliero
et al., 2019; de Lavenne et al., 2019; Choubin et al., 2019; Ecrepont et al., 2019; Zamoum
& Souag-Gamane, 2019; Prieto et al., 2019; Guo et al., 2021; Alipour & Kibler, 2018).
An underlying assumption behind regionalization is that similar basins have similar hydrologic behaviors. This implies that differences/similarities across catchments can be

classified into physical descriptors such as, climatology, geology, geomorphlogy, etc, with 98 the assumption that incorporating these descriptors will improve streamflow prediction. qq In other words, hydrological behaviors as predicted from models for different catchments 100 shall be based on similarities with regional information that is characterized by catch-101 ment physical descriptors. These approaches have been given a comprehensive review 102 in particular for PUB (Guo et al., 2021; Samaniego et al., 2017; Beck et al., 2016) and 103 can be grouped into model-dependent (process-driven) and model-independent (data-104 driven) methods, where 'model' denotes process-based models (Prieto et al., 2019). 105

106 Model-dependent methods give hydrologic predictions from process-based models. Information from the existing process-based hydrologic model is transferred to ungauged 107 catchments based on certain criteria that link gauged to ungauged catchments. In prac-108 tice, since those existing hydrologic models are calibrated to a specific catchment, this 109 relies on some strategy of information transfer. A typical application of a model-dependent 110 method implements a well-calibrated local hydrological process-based model and appro-111 priate connections among catchments. In the review paper by Guo et al. (2021), model-112 dependent methods can be classified into three categories: similarity based methods, re-113 gression based methods, and hydrological signature-based methods or some hybrid of each. 114

The model-independent approaches are data driven and do not rely on physical pro-115 cesses to simulate streamflow. Data driven methods learn how to predict streamflow from 116 weather drivers and catchment physical descriptors directly without involving any hy-117 drological process descriptions. Depending on either one or multiple catchments of data 118 used, the data driven model will learn localized or regionalized hydrologic behaviors re-119 spectively. A local model is referred to as the model using hydrologic data from only one 120 catchment. By contrast, when the hydrology data from multiple catchments are used and 121 those catchments cover a wide range of all available hydrologic behaviors, the model is 122 called a global model. 123

For data driven methods, one family is the neural network (Besaw et al., 2010; Hsu 124 et al., 1995). Besaw built an artificial neural network on one catchment and transferred 125 to another similar catchment without adaptation. It yielded unsatisfactory predictive 126 performance (Besaw et al., 2010). In recent years, the Long Short-Term Memory (LSTM) 127 networks (Hochreiter & Schmidhuber, 1997), one sub-family of neural networks, have shown 128 burgeoning applicability in streamflow prediction tasks (Kratzert et al., 2018). LSTM 129 based methods predict streamflow from antecedent weather drivers. Kratzert et al. (2019a) 130 have shown that using physical descriptors will train a universal global LSTM based model 131 that outperforms process-based individual models given the same forcing data. One of 132 the two versions of the LSTM developed by Kratzert et al. provides additional physi-133 cal interpretation, that is, basin similarities are preserved in the well trained machine 134 learning (ML) model. In gauged scenarios, Feng (Feng et al., 2020) embedded a global 135 LSTM within a data integration framework (using predicted discharge from previous day) 136 and found that it could marginally reduce prediction bias in regions with high flow auto-137 correlation. Frame showed that global LSTM outperforms the National Water Model (NWM) 138 (Frame et al., 2020). In the poorly gauged scenarios, Ma (Ma et al., 2021) showed that 139 fine tuning a global LSTM learned from data rich basins improved predictive performance 140 in poorly gauged basins in contrast to local models learned solely from limited data. 141

It bears emphasis that regionalization approaches, either model-dependent and model-142 independent, rely heavily on physical descriptors. However, to obtain a satisfactory re-143 gionalization performance, physical descriptors need to be sufficient such that process 144 complexities and associated scaling issues (Blöschl & Sivapalan, 1995) are encompassed. 145 146 Otherwise, catchment scale prediction will be handicapped by the lack of sufficient information. For instance, modeling hydrological behaviors at the small scale can be ac-147 complished by incorporating local processes with a few parameters. However, the incor-148 porated processes and parameterization need to be adjusted, either made simpler or more 149 complex, to model hydrologic behaviors at a larger scale. The same adjustment also oc-150

curs when modeling hydrological behaviors between global scale and local scale, upstream 151 and downstream. Accounting for these complexities and heterogeneities, sufficient phys-152 ical descriptors must be involved. For example, Drost and Mudersbach found that merely 153 incorporating landuse data with no additional physical descriptors provided little improve-154 ment to streamflow prediction and therefore may not benefit regionalization (Drost & 155 Mudersbach, 2021). However, due to the complexity of each catchment, such a complete 156 characterization to resolve hydrologic complexity is difficult and challenging (Beven, 2020). 157 This issue will be even more pronounced in applying models to data sparse regions where 158 physical descriptors are limited, or even unavailable. 159

#### $_{160}$ 3 Methods

161

#### 3.1 Long Short-Term Memory Network

Long short-term memory network (LSTM) (Hochreiter & Schmidhuber, 1997) is 162 a special type of recurrent neural network designed especially for modeling time series 163 predictions. Indeed, LSTM is the state-of-the-art deep learning model to predict stream-164 flow (Kratzert et al., 2018, 2019a; Frame et al., 2020; Feng et al., 2020; Ma et al., 2021). 165 In contrast to a traditional recurrent neural network, LSTM avoids gradient vanishing 166 or explosion (Bengio et al., 1994) and therefore preserves long term temporal dependen-167 cies for time series forecasting. This is achieved by using the gating architecture, which 168 explicitly controls information flow and updates system hidden features. This memoriz-169 ing mechanism and long term dependency allows LSTM to be well suited to model stream-170 flow on a catchment scale. In particular, weather inputs feed and alter catchment response 171 in various temporal scales. Although flooding season yields quick surface water response, 172 the streamflow in winter periods in northern climates tends to have much longer response 173 time because of involved snow and snowmelt processes. With the capability of the LSTM 174 to account for long term dependency, it automatically learns these streamflow behav-175 iors from data. Furthermore, it has been shown that some of the hidden features learned 176 by the LSTM resemble snow processes (Kratzert et al., 2018). 177

An LSTM maps a sequence of time series input into the response variable. In this paper, we consider an LSTM based architecture that uses input features  $(\mathbf{x})$  spanning T days to predict the observed discharge on the last day of the T-day window. The involved equations of an LSTM models are given below.

$$\boldsymbol{i}[t] = \sigma(\mathbf{W}_{\mathbf{i}}\boldsymbol{x}[t] + \mathbf{U}_{\mathbf{i}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{i}})$$
(1)

$$\boldsymbol{f}[t] = \sigma(\mathbf{W}_{\mathbf{f}}\boldsymbol{x}[t] + \mathbf{U}_{\mathbf{f}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{f}})$$
(2)

$$\boldsymbol{g}[t] = tanh(\mathbf{W}_{\mathbf{g}}\boldsymbol{x}[t] + \mathbf{U}_{\mathbf{g}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{g}})$$
(3)

$$\boldsymbol{o}[t] = \sigma(\mathbf{W}_{\mathbf{o}}\boldsymbol{x}[t] + \mathbf{U}_{\mathbf{o}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{o}})$$
(4)

$$\boldsymbol{c}[t] = \boldsymbol{f}[t] \odot \boldsymbol{c}[t-1] + \boldsymbol{i}[t] \odot \boldsymbol{g}[t]$$
(5)

$$\boldsymbol{h}\left[t\right] = \boldsymbol{o}\left[t\right] \odot tanh(\boldsymbol{c}\left[t\right]) \tag{6}$$

where  $\sigma(\cdot)$  is sigmoid function,  $tanh(\cdot)$  is the hyperbolic tangent function, and  $\bigcirc$  means 178 element wise multiplication. W, U, b are model parameters, which will be learned dur-179 ing optimization. Other variables in equations represent basic computation units involved 180 in the calculation. As gating variables, i[t], f[t], and o[t] are input gate, forget gate, and 181 output gate, respectively. They filter the information from the current and the previ-182 ous time stamp, then combine them to update cell state c[t]. c[t] underlines the intu-183 ition that motivates the LSTM design. c[t] is maintained serially and embeds the tem-184 poral contextual information, which is characterized in g[t], to then update the hidden 185 representation h[t]. The stacked input x enters the LSTM sequentially and alters the 186 information inherited from the previous timestamp. The previous information is stored 187 in cell states c[t] and hidden states h[t], both of which characterizes the system mem-188 ory. Cell states c[t] and hidden states h[t] are initialized as zero vectors and then grad-189

ually modified until the final date in T-day time windows is reached. After a linear trans-190 formation,  $\boldsymbol{x}[t]$  infuses with previous hidden state  $\boldsymbol{h}[t-1]$  and then is non-linearly trans-191 formed in i[t], f[t], g[t], and o[t] via a corresponding activation function. The previous 192 timestamp's cell state c[t-1] is updated with f[t] and then merges with an element-193 wise product of i[t] and g[t], which injects new information, to form a new cell state c[t]. 194 After another hyperbolic tangent activation, this new cell state c[t] merges with o[t] and 195 therefore updates the current hidden state h[t]. After the consecutive alteration of T time 196 stamps, the final hidden state h[T] is then transformed into the target variable, which 197 in our case is streamflow. 198

In the context of regionalization based streamflow prediction, both dynamic weather variables and static catchment physical descriptors as formulated in equation 7:

$$Q_t = f(\boldsymbol{x^d}, \boldsymbol{x^s}) \tag{7}$$

where  $Q_t$  is streamflow,  $x^d$  is weather input vector, and  $x^s$  is a d-dimensional vector of 199 physical descriptors. It bears emphasis that for a given catchment,  $x^s$  is assumed to be 200 temporally static, while  $x^d$  is temporally dynamic. We assume catchment physical de-201 scriptors do not vary with the time. In this paper, we consider two widely used LSTM 202 based models as illustrated in Figure 1. Namely, these two models are EA-LSTM and 203 CT-LSTM (Kratzert et al., 2019a), where 'EA' denotes entity awareness while 'CT' de-204 notes concatenation. These models differ in terms of how  $x^s$  is added into the network. 205 In CT-LSTM physical descriptors are added before LSTM cell, whereas in EA-LSTM, 206 they are used within the cell. For clarifications, the CT-LSTM refers to the normal LSTM 207 used in Kratzert et al.'s paper (2019a). We add prefix 'CT' to 'LSTM' to emphasize that 208  $x^s$  is concatenated with weather drivers before entering the LSTM cell.



Figure 1: LSTM family illustration. Figure is from "Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets" by Kratzert et al. (2019a), Hydrology and Earth System Sciences, 23, 5092 (Kratzert et al., 2019a)

209

210

#### 3.1.1 CT-LSTM

In CT-LSTM, at each timestamp, the dynamic weather input  $x^d$  is concatenated with the physical descriptors  $x^s$  to form the model input x[t]:

$$\boldsymbol{x}[t] = \begin{bmatrix} \boldsymbol{x}^{\boldsymbol{s}}, \boldsymbol{x}^{\boldsymbol{d}}[t] \end{bmatrix}$$
(8)

This model input enters the LSTM (equation 1 to 6), gets updated via the calculation of gates, and yields the final output - streamflow prediction. Through the calculation, physical descriptors are not placed within the LSTM cells or gates.

#### 3.1.2 EA-LSTM

First proposed in (Kratzert et al., 2019a), EA-LSTM (Entity Aware LSTM) uses a modified version of LSTM where input gate takes physical descriptors as input instead of input features as previously shown in Equation 1. The key idea here is to explicitly empower the LSTM to customize its learning ability for catchment-wise adaptation.

$$\boldsymbol{i} = \sigma(\mathbf{W}_{\mathbf{i}}\boldsymbol{x}^{s} + \mathbf{b}_{\mathbf{i}}) \tag{9}$$

$$\boldsymbol{f}[t] = \sigma(\mathbf{W}_{\mathbf{f}}\boldsymbol{x}^{\boldsymbol{d}}[t] + \mathbf{U}_{\mathbf{f}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{f}})$$
(10)

$$\boldsymbol{g}[t] = tanh(\mathbf{W}_{\mathbf{g}}\boldsymbol{x}^{\boldsymbol{d}}[t] + \mathbf{U}_{\mathbf{g}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{g}})$$
(11)

$$\boldsymbol{o}[t] = \sigma(\mathbf{W_o} \boldsymbol{x^d}[t] + \mathbf{U_o} \boldsymbol{h}[t-1] + \mathbf{b_o})$$
(12)

$$\boldsymbol{c}[t] = \boldsymbol{f}[t] \odot \boldsymbol{c}[t-1] + \boldsymbol{i}[t] \odot \boldsymbol{g}[t]$$
(13)

$$\boldsymbol{h}[t] = \boldsymbol{o}[t] \odot tanh(\boldsymbol{c}[t]) \tag{14}$$

As illustrated in Figure 1b and also equations 9 to 14,  $x^s$  enters the LSTM via in-215 put gates, learns customized embedding (equation 9) for each basin, and updates the cell 216 states recurrently at each timestamp. It therefore explicitly controls what modules in 217 LSTM respond to different catchments. This learned embedding will merge with other 218 gates (f[t], g[t], o[t]), whose alteration are contributed by only dynamic weather inputs 219  $x^d$ . This separated role of  $x^s$  and  $x^d$  in EA-LSTM splits the contributions towards stream-220 flow prediction from  $x^s$  in contrast to  $x^d$ . Additionally, the learned embedding affords 221 an opportunity to examine cross-catchment response in a global model, which was shown 222 to be close to the cross-catchment analysis using true basin characteristics (Kratzert et 223 al., 2019a). 224

#### **3.2** Data

Our experiments use the continental hydrology dataset, CAMELS (Catchment At-226 tributes and Meteorology for Large Sample studies) (Addor et al., 2017). The CAMELS 227 data set contains continuous meteorologic input, observed streamflow data, and catch-228 ment dependent spatially varying but temporally physical descriptors. CAMELS encom-229 passes a total of 671 watersheds across the contiguous US. Due to some watershed de-230 lineation errors (Addor et al., 2017), we followed the suggestion from Kratzert et al. (2019a) 231 to select 531 basins whose watershed boundaries are confirmed to be correctly delineated 232 without digital errors. Each watershed is supplied with observed discharge and climate 233 forcing data from remote sensing products (Daymet(Thornton et al., 2020), NLDAS(Xia 234 et al., 2012), MAURER(Maurer et al., 2002)), climate models, and data assimilation with 235 daily temporal resolution. Additionally, a corresponding hydrological model (SAC-SMA. 236 Sacramento Soil Moisture Accounting model) is well calibrated for each watershed and 237 its physical simulation is also available. Adopting such a wide distribution of watersheds, 238 CAMELS provides a comprehensive and detailed physical description of watersheds. Se-230 lecting only a subset of those features as suggested by Kratzert et al. (2019a), we choose 240 27 physical descriptors from climatology, geomorphology and geology perspectives to char-241 acterize and discriminate across watersheds (Table A1 in the Appendix A.). 242

These 27-d catchment physical descriptors are static vectors  $(x^s)$  characterizing 243 each catchment. We selected meteorological data from an updated version of MAURER 244 as model dynamic input  $(x^d)$ , which are daily precipitation, daily minimum air temper-245 ature, daily maximum air temperature, average short-wave radiation, and vapor pres-246 sure. The observed discharge from USGS is our target variable  $(Q^{O})$ . Both daily me-247 teorological weather inputs and discharge data cover a reasonably long record spanning 248 from 1980 to 2014. The data for each catchment was partitioned into training and test-249 ing periods while building and evaluating deep learning models. Some experiments in-250 volved using a subset of training years or a subset of basins, therefore, we specify a de-251 fault assessment scheme as to train a global model using data from 531 basins with 20 252

225

years of data. Under this assessment scheme, the training period starts from October
1st 1999 and ends on September 30th 2008. For a consistent evaluation, through all experiments, the testing period ranges between October 1st 1989 and September 30th 1999.

#### 3.3 General setup

256

Among different LSTM-based models, We apply the same optimization algorithm 257 (Adam optimizer (Kingma & Ba, 2017)) for training purposes to determine model pa-258 rameters. Model parameters are learned from data and are thus continuously updated 259 during training. The machine learning implementation also needs to specify hyper-parameters, 260 which are set before training without learning from data. During training, hyper-parameters 261 will not be updated. A few essential hyper-parameters include the look back period T262 and the dimension of hidden states h[t]. Adopting the previous work's specification (Kratzert 263 et al., 2019a) of these hyper-parameters, we determine T to be 270 days and the dimen-264 sion of hidden states to be 256. For the details on other hyper-parameters (e.g., learn-265 ing rate, batch size), please read the Appendix B in Kratzert et al.'s paper (2019a). 266

Machine learning models have uncertainties in model parameters after training. Ini-267 tialized randomly, model parameters will often be optimized to different values during 268 training. In simplistic terms, different model initializations will yield different models 269 after training. Accounting for uncertainty, it has been shown that ensemble results from 270 multiple model runs will facilitate the overall model performance (Kratzert et al., 2019a). 271 Therefore, the streamflow prediction result in all following sections is an ensemble mean 272 of five model realizations. For instance, the prediction of the EA-LSTM using physical 273 descriptors is an average of five model predictions, which are optimized from different 274 initializations. Note that for the Gaussian vector experiment, the randomness originates 275 from two sources, including model initializations and the Gaussian vector assignment. 276 For each of the five runs, their Gaussian vectors are assigned with different values. 277

Training deep learning models also requires a specification of the objective function. To account for cross-catchment variance, which is not considered in the commonly used mean squared error option, we use a smooth-joint NSE function (Kratzert et al., 2019a). The smooth-joint NSE function is shown below.

$$NSE^* = \frac{1}{B} \sum_{b=1}^{B} \sum_{t=1}^{N} \frac{(Q_t^m - Q_t^o)^2}{(s(b) + \epsilon)^2}$$
(15)

where *B* is the number of catchments, *N* is the number of daily data (days) for one catchment, which is indexed by *b*.  $Q_t^m$  is the predicted discharge at timestamp  $t(1 \le t \le N)$ , while  $Q_t^m$  is the corresponding observed discharge. s(b) is the standard deviation of the  $Q_t^o$  in basin *b* during training periods.  $\epsilon$  is a constant term ( $\epsilon = 0.1$ ) to avoid potential loss function explosion issue, which happens for catchments with extremely low s(b).

For consistent model comparison, we're using the NSE score instead of RMSE (root mean squared error) to evaluate streamflow prediction. NSE is a metric suited particularly to evaluate hydrological predictions.

$$NSE = 1 - \frac{\sum_{t=1}^{T} (Q_t^m - Q_t^o)^2}{\sum_{t=1}^{T} (Q_t^o - \bar{Q}^o)^2}$$
(16)

- $Q^m$  is predicted discharge,  $Q^o$  is observed discharge,  $\bar{Q^o}$  is the mean of observed discharge. A NSE score of 1 indicates a perfect time series prediction.
- 3.4 State of the Art

In terms of data-driven regionalization methods, CT-LSTM and EA-LSTM have been shown to perform satisfactorily for the streamflow prediction task (Kratzert et al.,

- 289 2019a). To remind readers of the state-of-the-art performance which relies on the phys-
- ical descriptors as shown in Table A1, in Figure 2, we show the testing NSE score for
- each catchment in the CAMELS dataset.

Table 1: State of the art LSTM based model. Mean and median refer to the summary statistics of the testing NSE scores across all 531 catchments in CAMELS.

Model	Mean	Median
Local LSTM	0.543	0.576
Global LSTM w/o static vectors	0.529	0.634
Global EA-lstm with 27-d descriptors	0.698	0.733



Figure 2: State of the art global regionalization performance using LSTM based deep learning architecture.

Local LSTM uses hydrologic data from only one catchment and does not need phys-292 ical descriptors  $(x^s)$  to combine data from multiple catchments. Thus, for 531 catchments, 293 there are 531 Local LSTM models. On the other hand, global LSTM refers to a global model learned from the training data of 531 catchments. While the Global LSTM merges 295 data from multiple catchments but does not use physical descriptors to adapt the net-296 work for different basins, the Global EA-LSTM with 27-d descriptors is also a global model 297 trained and tested using all 531 catchments but it takes advantage of 27-d physical de-298 scriptors to perform robust regionalization. As shown in Table 1, both the mean and me-299 dian of its NSE score is the highest (0.698 and 0.733 respectively) among the three model 300 options. In this gauged prediction scenario, cross-catchment information sharing ben-301 efits global training and thus elevates predictive performance. These results have been 302 previously shown by Kratzert et al. (2019a). 303

#### 3.5 Proposed Approach

304

In this paper, our aim is to answer the question "How to perform regionalization when catchment physical descriptors are unavailable, uncertain, or of insufficient dimension?" To address this issue, we propose to assign a vector of random values as a surrogate for missing physical descriptors. Since a set of random vectors doesn't have any similarity structure (i.e. correlation between any two random vectors is zero), they are a suitable baseline to incorporate the fact we don't have any prior information on catchment similarity due to missing physical descriptors. By using these random vectors, we enable the deep learning network to account for heterogeneity in catchment responseswhile sharing data across multiple basins.

Furthermore, the proposed concept of using random vectors as a baseline can also be used to evaluate the efficacy of known catchment physical descriptors. In other words, the performance difference between using random vectors and physical descriptors can imply the quality of physical descriptors. In section 4, we provide an extensive analysis of this concept in the context of streamflow prediction. In this paper, we consider two different strategies to create random vectors (Figure 3) as described below.

3.5.1 Gaussian Random Vectors

320



Figure 3: Random vector illustration. (a) refers to the d-dimensional Gaussian vector, while (b) illustrates the one-hot vector concept.

Figure 3a is a visual representation of the Gaussian vector (d-dimension) for all catchments. Random colors represent random numbers drawn from Gaussian distribution. In this strategy we assign d-dimensional vectors to each catchment where the vector values are drawn from a Gaussian distribution with zero mean and unit standard deviation. In other words, we randomly map each basin to a point in d-dimensional feature space.

326 3.5.2 One-hot Vectors

Figure 3b illustrates the one-hot vector representation. Each catchment is associ-327 ated with a binary vector that is 1 for one dimension and is zero elsewhere. The dimen-328 sion of the one-hot vectors equals the number of catchments. These one-hot vectors orig-329 inated from the binary vectors used to encode categorical variables in regression, where 330 in our case, the variable is catchment ID. There is one such one-hot binary vector for each 331 basin and these vectors are orthogonal to each other. It bears emphasis that there's no 332 freedom for the user to determine the dimension of the one-hot vector after the number 333 of catchments in a global model is known. For k basins, the length of the one-hot vec-334 tor for each basin is k. Although the one-hot vector does not involve random numbers, 335 the randomness in this random vector assignment is from basin order. Regardless of how 336 basins are sorted, one-hot vector assignment assures each basin will be assigned uniquely. 337

#### **4** Experiments and Results

We evaluate the effectiveness of our random vector approach with a series of experiments. First, in Section 4.1 we compare the random vector performance to that of the state-of-the-art EA-LSTM model. Next, we investigate the applicability of the random vector approach under varying data richness scenarios. In Section 4.2, we create a data inadequacy scenario by limiting the number of basins used in the training data. We also examine the impact of limiting the number of years of training data, as demonstrated by the experiment in Section 4.3. To further assess the generalizability of the random

vector approach, we evaluate other model settings in Section 4.4 and present our anal-346 vsis of the performance of the CT-LSTM model using random vectors along with the data 347 inadequacy scenario. Additionally, we compare the efficacy of the EA-LSTM and CT-348 LSTM models using random vectors. In Section 4.5, we explore the practical implica-349 tions of employing random vectors to model catchment complexities where physical de-350 scriptors for the system are incomplete. Finally, we show how the use of high-dimensional 351 representation of catchments improves regionalization by distinguishing them from one 352 another. 353

Implementing these experiments needs to specify a selective combination of the model 354 architecture (EA-LSTM or CT-LSTM) and static vectors  $(x^s)$ . Options for  $x^s$  include 355 27-d physical descriptors, random vectors, and mixing Gaussian vectors. For the sim-356 plicity of representing the results, we'll use acronyms to denote corresponding results of 357 those experiments, that is, the combination of model architecture and  $x^s$ . These acronyms 358 are shown in the table 2. Models for the incomplete physical systems are not given acronyms. 359 The comparisons across different models also involve statistical significance tests. Through-360 out the rest of this paper, we conducted pairwise statistical significance tests (Wilcoxon 361 signed-rank tests) to evaluate model results differences. The statistical significance is eval-362 uated using a 0.05 p-value threshold. 363

Table 2: This acronym table denotes the acronyms of model implementations. Combinations of model architecture and  $x^s$  specifications are shown in their acronyms. The "d" in these notations represent the dimension of  $x^s$ , which is only needed to specify the models using Gaussian vectors. For instance, EG-512 means EA-LSTM model using 512-d Gaussian vectors. "\*' means the corresponding models were not implemented.

$x^s$		EA-LSTM	CT-LSTM
27-d physical descriptors		EP	CP
Random	Gaussian d-dimension	EG-d	CG-d
vectors One-hot		EO	CO
Mixed Gaussian d-dimension vectors		EM-d	*

#### **4.1 Effectiveness of Random vectors**

To evaluate the applicability of our proposed random vectors method in regionalization, we first compared the predictive performance of a global model using random vectors (Gaussian or one-hot) against that using physically meaningful 27-d descriptors under EA-LSTM settings. The baseline model is the EP(Kratzert et al., 2019a) to show the state-of-the-art predictive performance. Substituting the 27-d physical descriptors with random vectors, our proposed method is implemented as either EG-d or EO.

Implementing Gaussian vectors requires a specification of *d*, which is determined empirically. The cumulative density function plot of the NSE score, shown in Figure 4 suggests using 512 (black solid line) as the Gaussian vector dimension because its testing performance is optimal compared to others.

The scatter plot (Figure 5) shows the testing NSE of the EG-512 and the EO versus the EP respectively across all 531 basins. Among these results, testing NSE scores less than -0.1 are forced to be -0.1 for illustration purposes. For each scatter plot, a cumulative density function (cdf) plot of NSE is also given. The EG-512 scatter plot is slightly upper skewed, the cdf of the EG-512 is also slightly right skewed compared to the EP. Figure 5a and Figure 5c shows that the EG-512 prediction performance is comparable



Figure 4: Cumulative density functions of the NSE score across different d Gaussian vectors for the EG-d. The X-axis is NSE score, which is truncated between 0.4 and 1 for a better illustration. The black dashed line represents the testing score corresponding to the EP. The black solid line corresponds to the EG-512, which yields the best performance in the EG-d.



Figure 5: Performance comparison cross the EP, EG-512 and EO. Model architecture is EA-LSTM. (a) shows the NSE score scatter plot that compares EG-512 and EP while its cdf comparison is shown in (c). (b) shows the comparison between EO and EP while its cdf comparison figure is (d).

381 382 to, if not slightly better than, the EP (statistically significant). In Table 3, the mean and median of the EG-512 is 0.711 and 0.746, both of which yield slightly more satisfactory results than the EP. The same comparison between the EO and the EP also yields a sim-

Table 3: Performance comparison of the EA-LSTM using random vectors against physical descriptors. Statistical summaries across all 531 basins are in column 'mean' and 'median'. The EG-512 and the EO are statistically different than the EP at 0.05 level.

Catchment static vectors	Mean	Median
27-d physical vectors (EP)	0.698	0.733
512-d random vectors (EG-512)	0.711	0.746
one-hot vectors (EO)	0.707	0.745

383

ilar trend. The mean and median of NSE score for the EO is 0.707 and 0.745. The EO

 $_{385}$  reaches comparable prediction performance to the EP (statistically different at 0.05 level).

For the statistical comparison of Table 3 and all following tables (Table 4 to Table 10), their corresponding p-value can be found in the supplemental information (Table S1 to

<sup>388</sup> Table S9).).

As we can see, using random vectors gives performance comparable to using known 389 physical descriptors. Furthermore, the random vector approach leads to significantly bet-390 ter results when compared to other strategies that do not use known physical descrip-391 tors (i.e. Figure 2, building local models or trivial merging of data from multiple basins) 392 Hence, the random vector approach is a viable solution when catchment characteristics 303 are not available. This performance is evaluated using the standard setting (section 3.2, 10 years training data from 531 basins). Although such abundant training data shows 395 slightly elevated testing performance, the proposed random vector method might still 396 be inapplicable in data poor situations. To assess the impact of data sparsity, we con-397 ducted an exhaustive analysis on the different data inadequacy scenarios with either fewer 398 number of basins or fewer number of training years. 399

400

#### 4.2 Effect of number of basins

For this situation we're creating a data inadequacy scenario where the training data consists of a limited set of k basins. Such a group of limited basins forms an insufficient global hydrologic dataset to train an LSTM based model. This experiment aims to answer the question "Given only k basins without physical descriptors, will the proposed random vector strategy be applicable for regionalization?" We vary k from 10 to 50 to 100 and follow the default assessment scheme as outlined in Section 3.2.

To generate the basin sets, we randomly select k basins as a group repetitively without replacement until all basins are selected. When the remaining basins cannot form a group with exactly the size k, those basins are either merged with the last group or form a stand-alone group as long as its order of magnitude approximates to k. For instance, when selecting 10-basin group, we select 53 groups in total, and the last group contains 11 basins. Similarly, the last group (11th group) in the 50-basin group has 31 basins. The last group (fifth group) in 100-basin group has 131 basins.

For the 53 groups of 10-basin groups, we compare the predictive performance us-414 ing random vectors relative to the performance of the model using 27-d physical descrip-415 tor. This comparison is illustrated in Figure 6. The X-axis denotes one-hot vector and 416 Gaussian vectors (varying d). Each category shows a box plot of performance compar-417 ison across basins. Median (blue dots), 25th percentile and 75th percentiles (upper and 418 lower box line) are shown for each box. Black hollow circles outside the upper and lower 419 box lines are outliers outside the specified quantile range. The Y-axis is the NSE score 420 improvement for each individual catchment compared to the 27-d physical descriptors. 421 The red line indicates the threshold for improved performance. A box plot whose NSE 422 distribution is skewed to positive NSE score improvement (above the threshold line) in-423 dicates a general performance improvement in that random vector category. Both the 424 EG-256 and the EG-512 show an performance improvement more pronounced than other 425 Gaussian vector dimensions and one-hot vectors. They both improve the NSE score on 426 an average of 0.082 (or in median 0.066). The reference performance of the EP is reported 427 at the table 4. The mean predicting NSE score is 0.308 while the median is 0.317, both 428 show that the 10-basin group downgrades the model performance because fewer basins 429 provide only limited training data and thus constrains model learning generalizable hy-430 drologic behavior. 431

For the 50-basin group and 100-basin group, the plot of NSE improvement is shown in Figure 7 except that we plot only the median of each case for a succinct visualization. The red line also marks the performance improvement threshold. Table 4 summarizes the NSE score improvement for all cases. As data from a greater number of basins are involved, the model performance gradually increases from 0.317 (10-basin) to 0.599 (50-



Figure 6: Random vectors implementation for 10-basin group EA-LSTM. Categories along the X-axis represent random vectors, including one-hot vectors (length of 10) and Gaussian vectors (dimension d varies from 2 to 1024). The Y-axis shows the individual basin NSE score improvement of the random vector in contrast to its corresponding EP, which is trained using the same basins. A zero NSE improvement indicates an improvement threshold marked by the red line. Within each category, 531 NSE improvement scores are distributed in the box plot where outliers exceeding 25th and 75th quantile are marked by black hollow circles.



Figure 7: Random vectors implementation for k-basin group EA-LSTM. k varies from 10 to 50 to 100. The median in Figure 6 are blue lines. Within each random vector category as shown in X-axis, the median of the individual-basin NSE score improvement in contrast to EP for k basins is plotted. Orange dots are the 50-basin group while green color represents 100-basin group

<sup>basin) to 0.656 (100-basin), all of which are lower than 0.733, the performance of the model
using all 531 basins. Note that both the trend and the performance are comparable to
the previous work, where the impact of the training data inadequacy on the EA-LSTM
performance is explored (Gauch et al., 2021). Table 4 and Figure 6 show a consistent
performance improvement comparison. Regardless of how limited the number of basins,
the Gaussian vector strategy (with an optimal dimension of either 256 or 512) slightly</sup> 

<sup>&</sup>lt;sup>443</sup> improves the 27-d physical vectors. In particular, the performance improvement from

the Gaussian vectors becomes saturated when d reaches 256 or 512. For 50-basin group,

Table 4: The individual-basin performance improvement of the EG-d and the EO to the EP. Note that the EP row does not show the NSE improvement, instead it shows the NSE score performance, which is the red dashed line performance in Figure 6 and 7. The largest performance improvement as indicated by the positive largest numbers is in bold font.

k-basin group		10	50	100	10	50	100
	d		mean		1	nedian	
	2	-0.073	-0.085	-0.09	-0.054	-0.076	-0.077
	8	-0.051	-0.061	-0.058	-0.033	-0.053	-0.046
	16	-0.032	-0.031	-0.03	-0.019	-0.029	-0.025
Constan	32	-0.015	-0.01	-0.008	-0.01	-0.007	-0.004
Gaussian	64	0.014	0.018	0.025	0.01	0.016	0.021
vector (EC 1)	128	0.047	0.045	0.039	0.035	0.036	0.031
(EG-d)	<b>256</b>	0.082	0.062	0.053	0.066	0.054	0.044
	512	0.082	0.055	0.053	0.066	0.048	0.046
	1024	0.06	0.038	0.039	0.04	0.036	0.034
one-h	not (EO)	-0.066	-0.035	-0.03	-0.043	-0.029	-0.022
$\overline{27-d \text{ physical}}$	descriptors (EP)	0.308	0.569	0.620	0.317	0.599	0.656

the average of the single-basin NSE improvement is 0.062 at 256-d while the median of 445 the single-basin NSE improvement is 0.54. For 100-basin group, EG-512 improves the 446 NSE score slightly with a mean of 0.053 while the median is 0.046, which is approximately 447 the same to the extent of what EG-256 improves. When the dimension of the Gaussian 448 vector becomes a higher 1024-d, the performance improvement begins to degrade as in-449 dicated by a smaller NSE improvement. In summary, we show that random vector ap-450 proach shows robust performance even with fewer number of catchments in the dataset 451 and hence can be used in situations where only few catchments are available. 452

#### 4.3 Effect of number of training years

453

In addition to the number of basins, another perspective on data inadequacy is the 454 number of training years. Varying the training years from 1 to 2 to 5 years, we sought 455 the answer to this question "Given only a few years of training data, will the proposed 456 random vector strategy be applicable for regionalization? " An LSTM model is trained 457 for all 531 basins with a limited number of years. The EG-d is tested against the EP un-458 der three sparse data cases, which are 1 year of data (October 1st 2007 to September 459 30th 2008), 2 years of data (October 1st 2006 to September 30th 2008) and 5 years of 460 data (October 1st 2003 to September 30th 2008). Models are tested for the same years 461 (October 1st 1989 to September 1st 1999) for consistent comparison. 462

Our previous empirical analysis indicates an optimal specification of d (Gaussian 463 vector dimension) to be 512 (Section 4.1), so the implementation of basin random vec-464 tors includes either 512-d Gaussian vectors or one-hot vectors. The Figure 8 shows that 465 both random vector strategies lead to prediction performance similar to the case utiliz-466 ing 27-d physical descriptors. For the reference, as the number of years of training data 467 increases, the performance of EP also increases from 0.632 (1 year) to 0.697 (two years) 468 to 0.766 (five years). This increasing trend was also identical to what Gauch et al. showed. 469 In particular, the EG-512 yields a more satisfactory performance than the EO. As shown 470 in Table 5, a NSE score improvement (both in mean and median) is observed when im-471 plementing 512-d Gaussian vectors, while the NSE score improvement is only observed 472 when using 5 years of training data when the one-hot vector strategy is applied. The re-473

sults show that even when training data are limited, randomly assigned vectors are still
able to learn as well as 27-d physical features.

Table 5: The impact of the number of training years on the performance improvement of random vectors for EA-LSTM. "Mean" and "Median" refer to statistics of the individualbasin NSE score improvement in relative to EP. The EP row does not show NSE score improvement, instead it shows the NSE score performance, which is the reference performance in Figure 8. Positive numbers mean that random vectors yield better predictive performance.

Number of training year	S	1	2	5
Gaussian 512-d (EG-512) one-hot (EO)	mean median mean	$\begin{array}{r} 0.013 \\ 0.009 \\ -0.025 \\ 0.022 \end{array}$	$\begin{array}{r} 0.052 \\ 0.041 \\ -0.003 \\ 0.005 \end{array}$	$\begin{array}{r} 0.026 \\ 0.019 \\ 0.015 \\ 0.012 \end{array}$
27-d physical descriptors (EP)	mean median	0.399 0.632	0.628 0.697	0.013 0.719 0.766



Figure 8: The impact of the number of training years on EA-LSTM. The Y-axis represents the individual-basin NSE score difference between the corresponding category in X-axis and the predictive performance using 27-d physical descriptors (EP). The red line indicates performance improvement threshold.

## 4.4 Performance of alternative models

476

As outlined in the section 3.1, both EA-LSTM and CT-LSTM adopt  $x^s$  in different ways. From previous sections (section 4.1, 4.2, and 4.3), we've shown the efficacy of the random vectors in EA-LSTM in both data rich and data poor scenarios. It remains unknown whether random vectors are applicable to the CT-LSTM. Experiments of this section are designed to clarify this doubt.

We first evaluated the performance of random vectors under the CT-LSTM setting to answer this question "Under different model architectures, will the proposed random vector strategies be applicable for regionalization?". We then examined the regionalization performance of random vectors across the EA-LSTM and the CT-LSTM to answer the question "Which random vector strategy is better suited for regionalization, Gaussian vectors or one-hot vectors?" Last, we selected the CT-LSTM as the model architecture for an exhaustive analysis on the data inadequacy cases in terms of basin numbers.

#### 4.4.1 Random vectors in the CT-LSTM

490

For the CT-LSTM, the Gaussian vector implementation needs to specify the op-491 timal vector dimension d. Figure 9 shows that the CG-16 yields the most satisfactory 492 performance among different Gaussian vector dimension options. Therefore, we empir-493 ically select 16 as the optimal Gaussian dimension to represent the CG-d performance 494 (Figure 10c). Note that the optimal 16-d of the CG-d is less than the optimal 512-d of 495 the EG-d. We'll explain this in the section 5.1 in "Discussion" section. Using 27-d phys-496 ical descriptors, CP achieves performance comparable and slightly better than EP (Fig-497 ure 10a and Table 6). The median NSE score performance improves from 0.733 (CP) to 498 0.744 (CO). Random vector options (CO and CG-16) slightly outperform 27-d physical 499 descriptors (CP). The median of testing NSE performance improves from 0.744 to 0.754 500 when using the one-hot vector strategy, while the CG-16 elevates the performance to 0.752. 501



Figure 9: Cumulative density function plots of the NSE score across different d Gaussian vectors for the CT-LSTM. The X-axis is truncated between 0.4 and 1 for a better illustration. The black dashed line represents the testing score of the CP, the black solid line corresponds to the optimal 16-d performance among the Gaussian vector groups (CG-16).

Table 6: Random vector comparison cross different models (The CT-LSTM based random vector performance is statistically different from the CP at 0.05 level).

models	Mean	Median
EP	0.698	0.733
CP CO	$0.715 \\ 0.720$	$0.744 \\ 0.754$
CG-16	0.717	0.752



Figure 10: Predicted performance comparison of a random vector implementation in CT-LSTM (CO and CG-16) in contrast to CT-LSTM using 27-d physical descriptors (CP). (a) is the comparison between CP and EP; (b) is the comparison between CO and CP; (c) is the comparison between CG-16 and CP.

Although the slight improvement of the CO and the CG-16 in contrast to the CP imply the applicability of random vectors in the CT-LSTM, data abundance has always been an important factor impacting the machine learning model performance. To consolidate the argument that CT-LSTM with random vectors, especially one-hot vectors, yields better performance consistently under various data richness scenarios, we repeated the experiments outlined in section 4.2 for the CT-LSTM. Training data are limited by the number of basins.

Table 7: The improvement of random vectors over 27-d physical features in the CT-LSTM. "Mean" and "Median" refer to statistics of NSE score improvement for individual basins in relative to the CP. Note that the CP row shows the NSE value for the CP, instead it shows the NSE score performance, which is the reference performance in Figure 11 and 12. The most satisfactory performance is in bold font: 32-d Gaussian vector, 64-d Gaussian vector, and one-hot vector. For 10-basin and 50-basin group, their NSE performance difference between random vectors and the CP counterpart is significantly different at 0.05 level, while the NSE difference comparison at 100-basin group does not show statistical significance.

k-b	asin group	10	50	100	10	50	100
	d		mean		r	nedian	
	2	-0.079	-0.073	-0.074	-0.075	-0.063	-0.063
	8	-0.055	-0.031	-0.024	-0.050	-0.028	-0.023
	16	-0.037	-0.015	-0.007	-0.037	-0.018	-0.010
Caucion	32	-0.022	-0.006	0.004	-0.023	-0.008	0.001
Gaussian	64	-0.023	-0.004	0.002	-0.021	-0.006	0.000
(CC - 1)	128	-0.041	-0.019	-0.019	-0.039	-0.016	-0.003
(UG-d)	256	-0.074	-0.059	-0.033	-0.072	-0.048	-0.026
	512	-0.103	-0.125	-0.094	-0.097	-0.114	-0.083
	1024	-0.134	-0.188	-0.174	-0.129	-0.191	-0.171
one	-hot (CO)	-0.046	-0.007	-0.005	-0.042	-0.005	0.001
27-d physic	al descriptors (CP)	0.454	0.655	0.684	0.481	0.687	0.709

509 510 511 Figure 11 exhibits a box plot showing the NSE improvement for the 10-basin group using the CT-LSTM architecture. Any point above the red line (NSE score improvement threshold) indicates a performance improvement in contrast to 27-d physical descriptors.



Figure 11: Impacts of random vectors on CT-LSTM for a 10-basin group. Categories on the X-axis represent random vectors, including one-hot vectors (length of 10) and Gaussian vectors (dimension d varies from 2 to 1024). The Y-axis show the NSE score improvement for individual basin of the random vectors in contrast to the CP. A zero NSE improvement indicates no performance improvement marked by the red line.

- In the 10-basin group category, the optimal Gaussian d for the CG-d is lower than that
- of the EG-d. The optimal Gaussian vectors performance is comparable to that of one-
- <sup>514</sup> hot vectors. To obtain a general insight, we varied k from 10 to 50 to 100 and therefore produced the following result in Figure 12 and Table 7.



Figure 12: A random vector implementation for a k-basin group CT-LSTM. k varies from 10 to 50 to 100. Within each random vector category as shown in X-axis, the median of NSE improvement score for individual basin in contrast to the CP for k basins is plotted. Blue dots are 10-basin group, orange dots are 50-basin group, while green color represents 100-basin group

515

Figure 12 shows the median of the NSE improvement using random vectors in the 516 CT-LSTM in contrast to the CP. Dots below the red line mean the prediction perfor-517 mance of the corresponding categories is worse than the CP. As the number of catch-518 ments available for training increases, the one-hot vector strategy and optimal Gaussian 519 vectors in the CT-LSTM yields performance comparable to the CP. The optimal d for 520 the CG-d is either 32 or 64, which is lower than the optimal 512-d in the EG-d. As also 521 recognized in Figure 10, this discrepancy of optimal Gaussian d between the CT-LSTM 522 and EA-LSTM can be explained by the number of parameters involved in these model 523 architectures and we'll expand this discussion in section 5.1. We point out that these ran-524 dom vector strategies are approximate to but do not marginally exceed the CP perfor-525 mance. In particular, the relative significant performance improvement occurs when us-526

ing the one-hot vector in the 100-basin group but to a much lesser extent. Varying k from
10 to 50 to 100, as more catchments are involved until 531 basins are included, the onehot vector is a preferable random vector strategy for CT-LSTM than Gaussian vectors.

530

547

#### 4.4.2 Best performance of using random vectors

In the CT-LSTM setting, the above experiments demonstrate that random vector strategies still prevail over 27-d physical vectors. The best random vector strategy for the CT-LSTM is one-hot vector, while the best random vector method for the EA-LSTM is 512-d Gaussian vectors. The preferable random vector strategy varies depending on the model setting.

In a pursuit of model performance when utilizing random vectors, we need to pro-536 vide a practical solution to the question "When implementing random vectors to per-537 form regionalization, shall I use Gaussian vectors or one-hot vectors?". We next com-538 pared the optimal random vector performance between the CT-LSTM and EA-LSTM. 530 Figure 13 shows the testing NSE difference from the various EG-d against the CO. Based 540 on the previous result showing that the EO is not as good as EG-d, 'one-hot' on X-axis 541 (EA-LSTM random vector strategy) is omitted. Figure 13 shows the median of the NSE 542 difference for various selections of k basins. All points are below the performance thresh-543 old line, indicating that the CO slightly outperforms EG-d. When implementing the ran-544 dom vector strategy as a surrogate for missing physical descriptors, the best performance 545 is obtained when applying CO. 546



Figure 13: The performance of EG-d in contrast to the CO for k-basin group. k varies from 10 to 50 to 100. The Y-axis is the NSE difference score quantifying the performance of the EG-d (categories in X-axis) relative to the CO for individual basin. The red line marks no NSE difference. Plotted points are median of the NSE difference across all basins. For points below the red line, they mean that the CO yields more satisfactory performance

#### 4.5 Incompleteness of physical descriptors

So far we considered the scenario where physical descriptors are not available and 548 assessed the performance of our random vector approach. In this section, we consider 549 a more common regionalization challenge where physical descriptors are incomplete. In-550 complete physical descriptors under-represent a system of catchments and can only help 551 regionalization in a limited degree. To tackle this problem, the question becomes "will 552 the proposed random vector strategy benefit the model regionalization in this informa-553 tion deficient physical system?" To assess the performance caused by this deficiency, we 554 define a physically underrepresented global system in CAMELS where only a subset of 555 27-d physical descriptors is used to distinguish basins. We compare the global LSTM us-556

ing random vectors in contrast to the global model using these insufficiently informative
 descriptors. One extreme case is a system without any static catchment descriptors, which
 has been shown in the section 3.4 (Figure 2).

Ignoring the model selection differences, we select EA-LSTM for this experiment because it explicitly modulates LSTM via static vectors. The EA-LSTM using some subset of 27-d physical descriptors is trained and compared. EA-LSTM using 9-d climate features, 10-d geology features, and 8-d geomorphology features are trained separately and compared to the EA-LSTM using random vectors.

Catchment hydrologic models are formulated to resolve complexity and associated 565 scaling issues in hydrological processes. Both issues will require a comprehensive phys-566 ical understanding. From a practical perspective, static physical descriptors (for instance, 567 Table A1) can only characterize complex catchments to a limited dimension because a 568 sufficient catchment complexity characterization is challenging across scales. In the field scale, a hydrological model might characterize local hydrological processes completely, 570 but the applicability of this locally built model to a larger basin might fail if the model 571 is not adjusted, either simplified (reduce the number of parameters) or made more com-572 plex (enrich physical parameters). Therefore, for the regionalization involving catchments 573 at various scales, the question becomes "Are any given physical descriptors sufficient for 574 modeling the complexity of catchments?" This question also implies another question: 575 "how many physical dimensions do we need for characterizing the complexities of stream-576 flow generation processes?"



Figure 14: The performance of the EG-512 in contrast to EA-LSTM under a physically incomplete catchment system. Scatter plots show testing NSE scores comparison between the Y-axis and the X-axis. The y label is EG-512 and is fixed in the above 4 figures. The X-axis changes from a no physical feature system, a climate feature system, a geology feature system, to a geo-morphology feature system. The below cumulative density function plot collects NSE scores together for each category aforementioned. The black line is EG-512. We also plotted the benchmark performance with sufficient 27-d physical descriptors (grey solid line, EP) to remind the reader of the performance under a physically sufficient catchment system.

Table 8: Physical system completeness identification. The most satisfactory performance is in bold font (EG-512). In contrast to the EP, EG-512 and other physically incomplete system shows statistically different performance (at 0.05 level).

catchment static vectors	Mean	Median
(0-d) Without static features	0.529	0.634
512-d Gaussian vectors (EG-512)	0.707	0.745
9-d climate features	0.611	0.665
10-d geology features	0.638	0.679
8-d geomorphology features	0.630	0.680
27-d physical descriptors (EP)	0.698	0.733

To answer these questions, we compared our random vector results to models uti-578 lizing incomplete sets of physical vectors. In Table A1, 27-d physical descriptors are cat-579 egorized into three groups: climate, geology and geo-morphology. Among these, the de-580 scriptors of any single group are an under-representative description for basins. For in-581 stance, 9-d climate descriptors presumably characterize basins less informatively than 582 27-d physical descriptors. For this experiment, we choose the EA-LSTM as the model 583 structure and use 512-d Gaussian vector as its optimal random vector strategy. Each one 584 of the three descriptor subset groups leads to an EA-LSTM under a physically uninfor-585 mative system since complexities are simplified and the system incurs information loss. 586 For the extreme case where there are no physical descriptors present, the global model 587 is a simple global LSTM without basin characteristics, results of which were shown early 588 in section 3.4 (Figure 2). 589

In Figure 14, a distribution of scatters above the diagonal line (exactly equal per-590 formance from the methods indicated by axes) indicates that Gaussian 512-d vectors out-591 perform all these physically incomplete conditions. This fact is better illustrated in the 592 cumulative density function plot as the distribution of NSE scores is skewed to upper 593 tail. Both its mean and median NSE scores are higher than any physically incomplete 594 characterization (Table 8). Note that as shown earlier, the EG-512 case reaches compa-595 rable and slightly better performance than EP. This observation also implies that 27-596 d physical vectors are lacking additional physical characterizations. 597

598

## 4.6 Effectiveness of Distinguishing Basins in the High-dimensional Space

Missing or incomplete physical descriptors make catchments less distinguishable from each other. Even for the assumed complete 27-d physical descriptor, they suffer from losing information as heterogeneous catchment systems are spatially simplified. Of the 27-d physical descriptors, the spatially dependent descriptors are deterministic representations of catchments, such as soil porosity, silt fraction, etc. This simplification also reduces the functionality of static vectors to distinguish catchments from each other, which might produce disadvantages for regionalization.

Further, the proposed random vector strategy projects catchments in the high-dimensional space. In particular, the EG-512 assigns a Gaussian vector of 512-dimension to catchments, while the CO uses the 531-d one-hot encoded vector to represent catchments. In other words, it seems that characterizing catchments in a high-dimensional space distinguishes them from each other and thus improves regionalization. Recognizing this, the question becomes: "Can we incorporate 27-d physical descriptors in the high-dimensional space? "

We offer two methods to include the 27-d physical descriptors in the high dimen-613 sional space and compare the performance of these methods with the performance of the 614 random vector approach. We use EA-LSTM in these experiments rather than CT-LSTM 615 because training CT-LSTM increases the computational burden and complicates machine 616 learning. Additionally, EA-LSTM explicitly modulates the LSTM architecture. In the 617 first method, we concatenate the 27-d physical descriptors with additional Gaussian vec-618 tors to expand the dimension of  $x^s$ . We refer to this as the mixed Gaussian vector ap-619 proach. In the second method, we create an embedding layer to explicitly project  $x^s$  into 620 a high dimensional space before entering the EA-LSTM cell. 621

#### 4.6.1 Mixed Gaussian Vector

The 27-d physical vectors are augmented with extra dimensions filled by Gaussian 623 vectors, which is named as "mixed Gaussian vectors" (denoted as 'EM'). Catchments 624 are gradually more distinguishable as their Gaussian vector dimension increases. These 625 appended Gaussian vectors improve the distinctiveness of catchment characterization. 626 We define a global system with 64-d, 128-d, 256-d, and 512-d vectors, all of which in-627 clude the 27-d physical descriptors. For instance, for the 64-d  $x^s$ , besides the 27-d phys-628 ical descriptors, 37-d (64-27=37) vectors are randomly drawn from the Gaussian distri-629 630 bution.



Figure 15: Comparison of the performance between Gaussian vectors (EG-d, blue box) and mixed Gaussian vectors (EM-d, flaxen box). The X-axis is the dimension of the static vector (from 64 to 1024), while the Y-axis shows the NSE difference in contrast to the EP for individual basin. The red line specifies the performance improvement threshold. Box portions above the red line indicate performance improvement.

As shown in Figure 15, compared to the baseline performance, which is the EP, the mixed Gaussian vector (EM-d) yields better performance and achieves the maximal performance improvement at EM-512. On average, the NSE improvement is 0.018. From 64-d, to 1024-d, all of the EM-d results yield better performance (positive NSE score improvement statistics in Table 9). In contrast to a pure random Gaussian system, given the same *d* Gaussian dimension (blue and flaxen box in the same X-axis category), the Table 9: NSE performance difference of the mixed Gaussian vectors (EM-d) and Gaussian vectors (EG-d) in contrast to 27-d physical vectors (EP) for individual basin. Positive scores mean that the EP yielded worse predictive performance (these results are statistically different from the EP at 0.05 level). The highest value is highlighted in bold.

static vector dimension (d	l)	64	128	256	512	1024
Gaussian vectors (EG-d)	mean median	-0.007 -0.006	$0.010 \\ 0.002$	$\begin{array}{c} 0.010\\ 0.004\end{array}$	$\begin{array}{c} 0.009 \\ 0.007 \end{array}$	0.002
Mixed Gaussian vectors (EM-d)	mean median	$0.009 \\ 0.004$	$0.013 \\ 0.006$	$0.014 \\ 0.007$	$\begin{array}{c} 0.018\\ 0.011\end{array}$	$0.009 \\ 0.008$

EM-d also marginally improves the NSE score. In Table 9, the NSE improvement in 'Mixed
Gaussian vectors' are consistently more pronounced than 'Gaussian vectors' across varying static vector dimension d. The mixed Gaussian vector leads to marginally better global
model performance compared to either pure random Gaussian vector system or pure physical system. As such, it suggests that when physical descriptors are augmented with Gaussian
sian vector in a high dimension space, catchments are more distinctively represented, which
supports and benefits regionalization.

#### 4.6.2 Additional Embedding

644

Besides the mixed Gaussian vector approach, the other method to characterize catchments in a high dimensional space is to create an embedding layer between  $x^s$  and the input gate *i* of the LSTM cell. That is to say, the equation 9 in EA-LSTM is replaced by 17 and 18:

$$\boldsymbol{v} = \sigma(\mathbf{W}_{\mathbf{v}}\boldsymbol{x}^{\boldsymbol{s}} + \mathbf{b}_{\mathbf{v}}) \tag{17}$$

$$\boldsymbol{i} = \sigma(\mathbf{W}_{\mathbf{i}}\boldsymbol{v} + \mathbf{b}_{\mathbf{i}}) \tag{18}$$

where v is the embedding layer that  $x^s$  is mapped into. We chose 512 as the dimension of v for its empirical outstanding performance in the EG-512. We first mapped the physical descriptors into a 512-d embedding layer (denoted as 'PEA'). Then as its random counterpart when  $x^s$  is not available, we also experimented to map random vector into the embedding layer (denoted as 'REA'). To preserve the modeling capacity without including additional model parameters, the dimension of the random vectors in the REA is still 27.

Table 10: Performance comparison of the EA-LSTM using additional embedding layers against previous random vector approaches. Statistical summaries across all 531 basins are in column 'mean' and 'median'. This result is statistically different from the EP at 0.05 level.

Catchment static vectors	Mean	Median
27-d physical descriptors (EP)	0.698	0.733
512-d Gaussian vectors (EG-512)	0.711	0.746
27-d physical descriptors with 512-d embedding (PEA)	0.713	0.753
27-d Gaussian vectors with 512-d embedding (REA)	0.714	0.757

As shown in Table 10, the mean of testing NSE scores for the 'PEA' is 0.713 while its median is 0.753. With similar performance, the mean and the median of the testing NSE scores for the 'REA' is 0.714 and 0.757 respectively. Note that both 'PEA' and 'REA' shows slightly better performance than EG-512. In particular, 'PEA' and 'REA' yields comparable performance.

## 557 5 Discussion

669

This section is organized as four subsections. The first section (5.1) presents com-658 parison across the experiments in Section 4.1, 4.2, 4.3, 4.4 and shows the presence of an 659 optimal large d. The second section (5.2) compares the results in Section 4.5 and 4.6 and 660 discusses the regionalization advantage from high-dimensional characterization of catch-661 ments. The third section (5.3) presents the analysis of the embedding layers in both EG-512 and EP. It also shows the discussion on what all of our results will imply to under-663 standing catchment similarities and complexities. The forth section (5.4) presents a dis-664 cussion on the impact of our results towards the understanding of deep learning in the 665 context of streamflow prediction, and shows one practical utility of random vectors in 666 assessing the completeness of physical descriptors. The fifth section (5.5) outlines the 667 limitations of the current study and describes possible future directions. 668

5.1 Random Vectors

Results in section 4.1 show that the proposed random vector method achieves a performance comparable to the state-of-the-art model (Figure 5, Table 3). In other words, without any knowledge of physical descriptors, the global LSTM based model using random vectors successfully learns universal hydrologic behavior and sustains benchmark streamflow prediction performance. These random vectors retain practical feasibility without having to obtain any physical descriptions of basins. This is arguably the most significant scientific contribution of this paper.

The exhaustive analysis from section 4.2, 4.3, and part of 4.4 verifies the applicability of employing random vectors in data scarce regions. For a limited number of basins (Figure 6, 7, 11, 12, Table 4, 7) and a few years of training data (Figure 8), the two situations which restrict hydrologic extrapolation across catchments, random vectors are still viable for hydrologic regionalization.

The prediction performance achieved by random vectors varies between the EA-682 LSTM and the CT-LSTM, which we hypothesize is the result of different modulation 683 levels in their architectures. Random vectors functionalize as static vectors  $(x^s)$  that represent each basin, from which the LSTM family models modulate its internal computa-685 tion and mapping across neurons for each basin distinctively. That is, for a given weather 686 input  $x^d$ , the global model is aware of which basin the  $x^d$  data originates from and thus 687 modulates how streamflow shall be predicted differently in contrast to other basins. Yet 688 this modulation extent likely varies between EA-LSTM and CT-LSTM. CT-LSTM con-689 catenates  $x^s$  with  $x^d$  at each timestamp and thus performs a stronger modulation be-690 cause this catchment awareness is passed through all gates in the LSTM. Merely feed-691 ing into the input gate, the EA-LSTM does not modulate the network as well as the 692 CT-LSTM but it ensures  $x^s$  is involved in the temporal context update (memorizing and 693 forgetting). Across those different LSTM model selections with various catchment mod-694 ulation degree, random vectors consistently perform as well as, if not better than, phys-695 ical descriptors for learning across basins (Table 3, 6). 696

Random vectors can either be Gaussian vector or one-hot vector. For the Gaussian vector, note that its implementation needs to specify its dimension d, which is empirically obtained. The optimal d for the Gaussian vector is different between the EA-LSTM and the CT-LSTM. For the EA-LSTM, the optimal d is either 256 or 512 (Fig-

ure 4, 7, 6), while for the CT-LSTM, it is in the range of 16 to 64 (Figure 9, 11, 12). This 701 means that the EA-LSTM needs a higher dimension of static vectors to perform region-702 alization than the CT-LSTM. We explain this difference by the amount of trainable ma-703 chine learning parameters. The increasing number of trainable model parameters of CT-704 LSTM hinders the training processes. For the CT-LSTM, an increased static input will 705 expand the concatenated input  $\boldsymbol{x}[t]$  dimension (Equation 8), which in turn enlarges the 706 dimension of the transformation matrices  $W_i, W_f, W_g, W_o$ . In contrast, the static in-707 put  $(x^s)$  dimension only impacts EA-LSTM's input gate dimension  $(\mathbf{W}_i)$ . Consequently, 708 given the same  $x^s$  dimension augmentation, the parameter increment of CT-LSTM is 709 four times the increase of the number of parameters in EA-LSTM. A higher d-dimension 710 Gaussian vector CT-LSTM becomes more difficult to optimize than that for the EA-LSTM. 711

Although the optimal d differs between the CT-LSTM and the EA-LSTM, the per-712 formance saturation trend is identical. As illustrated by the results between section 4.1 713 and 4.4, when expanding static vector dimension, the predictive performance saturates 714 at a certain point and then deteriorates. This pattern indicates a presence of the opti-715 mal d, which cannot be too large or too small and is suitable for achieving best model 716 performance. In particular, the optimal d is universal regardless of the number of basins 717 involved (Figure 6, 7, 11, 12) as well as the number of training years (Figure 8). Thus, 718 it suggests implications for addressing catchment modeling complexities, which are of-719 ten entangled with associated scaling issues between catchments as one of the Two Clouds 720 in hydrology (Beven, 1987). Traditional hydrological models need to either be simpli-721 fied or made more complex to account for scaling transformations between catchments 722 that have different complexities. This can be done by reducing or increasing the num-723 ber of hydrologic parameters, which can also be reasonably interpreted as the dimension-724 ality of static vectors. A recognized optimal d in deep learning models illustrates that 725 the level of an appropriate scale for regionalization exists and the involved cross-catchment 726 hydrologic complexities exceed what the physical descriptors can provide. 727

The performance of the CT-LSTM with one-hot vectors (CO) is slightly better than 728 the CT-LSTM with Gaussian vectors (CG-16) as shown in Table 3, 7. To the EA-LSTM, 729 Gaussian vector is a more suitable choice yielding better performance. Although the spe-730 cific random vector strategy is favored by different LSTM choices, the outperformed ran-731 dom vector strategy all enjoys an advantage of high dimension characterization. On a 732 continental scale for 531 catchments, the one-hot vector is a vector with a length of 531, 733 while the optimal-d is 512, both of which far exceed the 27-d physical descriptors. The 734 high-dimensional static vector  $x^s$  enhance the LSTM's ability to learn across basins to 735 a similar or even better extent than what the 27-d catchment physical descriptors are 736 capable of performing. This insight and discovery has a broad and significant implica-737 tion for hydrologists to examine the value of  $x^s$  (either physical descriptors or random 738 vectors) that were brought in for modeling catchment complexities. Specifically, a more 739 relevant hydrologic question to ask is: "Are catchment physical descriptors sufficient to 740 model streamflow generation complexities and gauged catchment systems? If not, how 741 many dimensions do we need?" 742

743

#### 5.2 High Dimensional Catchment Characterization for Regionalization

The performance using 27-d physical descriptors is slightly worse than that of the 744 512-d random vectors. We interpret this to mean that the high dimensional character-745 ization of catchments benefits regionalization performance. This idea is further supported 746 by the result in Section 4.5, where we see a certain subset of 27-d physical descriptors 747 is also outperformed by 512-d Gaussian vectors. When physical descriptors are incom-748 plete, catchments are less distinguishable from each other and the regionalization per-749 formance worsens, as shown in Figure 14. Also note that Kratzert et al. (2019a) pointed 750 out that the 27-d physical features utilized in their study are intrinsically uncertain since 751 spatial heterogeneities are simplified as spatial averages and therefore lose certain regional 752

information. Uncertainties in physical descriptors might be another source that down grades the distinctiveness across basins.

Therefore, to create a system where a system of catchments can be more distinguishable from each other, or the complexity of catchment systems can be more sufficiently quantified, we showed two strategies to expand the dimension of physical descriptors, the mixed Gaussian vector (Section 4.6.1) and the use of additional embedding (Section 4.6.2).

The mixed Gaussian vector concatenates physical descriptors with additional Gaus-760 sian vectors. The added Gaussian vectors do not have any physical meaning and only 761 fill the expanded dimension with a Gaussian random number. Overall, with physical de-762 scriptors, this high-dimensional mixed static vector preserves the physical hydrological 763 information and the randomness simultaneously. Results (Figure 15, Table 9) in section 764 4.6.1 show that the high-dimensional mixed Gaussian vectors effectively distinguish catch-765 ments and thus improve the regionalization. The peak performance is realized by 512-766 d, which shows the largest NSE score improvement. The results also indicate that mixed 767 Gaussian vectors always outperform pure Gaussian vectors. Given the same dimension 768 of static vectors, the information contained in 27-d physical features improves regional 769 modeling. In contrast to a pure random system formed by all dimensions of Gaussian 770 vectors, the mixed Gaussian vectors introduce ordered information and physical simi-771 larities, and thus benefit regionalization. 772

Inserting an embedding layer before the input gates allows an opportunity to learn 773 more information in  $x^s$  as  $x^s$  is transformed into a high-dimensional space (512-d). When 774  $x^s$  is physical descriptors (PEA), its performance is approximate to the EG-512 and also 775 much better than the EP. After the transformation, the information of the original 27-776 d physical descriptors is extracted in a way that benefits regionalization. Meanwhile, this 777 performance improvement is also recognized to the case where  $x^s$  is 27-d Gaussian vec-778 tors (REA). Note that the modeling capacity for both PEA and REA is identical since 779 both have exactly the same amount of training parameters. This observation illustrates 780 that, within a context of the gauged basin scenario, a higher dimensional space with op-781 timal d-dimension where catchments become more distinguishable will always elevate the 782 regionalization performance regardless of how such a high dimension space originates (ei-783 ther from physical descriptor or Gaussian vector). 784

#### 785

## 5.3 Implications for Catchment Systems

Besides the discussion on evaluating the effectiveness in random vectors and the
high dimensional characterization advantage, the random vector approach itself and its
mapping mechanisms has significant implications for understanding the current work of
modeling hydrologic regionalization using the LSTM based models.

Both the Gaussian vectors and the one-hot vectors map catchments into a high di-790 mensional space. The Gaussian vector characterizes catchments as statistically indepen-791 dent from each other. For the one-hot vectors, the static vectors of catchments are or-792 thogonal to each other. Although these random vectors do not quantify catchment sim-793 ilarities, they ensure catchments are different from each other in a consistent way. This 794 implies that distinguishing catchments in a gauged system achieves state-of-the-art re-795 gional modeling performance in a deep learning framework, which reflects a recently raised 796 hydrologic concern expressed by Beven (2020) "When essential catchment characteris-797 tics are not well understood or defined and thus not even included in catchment phys-798 799 ical descriptors, how could a derived deep learning model perform satisfactory regional*ization performance?*" Although not explicitly defining catchment characteristics, our 800 proposed random vectors can be interpreted as high-dimensional non-physical descrip-801 tors characterizing the complexity of catchments. Catchment systems are composed of 802 linked components representing the functional relationships between weather inputs and 803



Figure 16: Clustering maps for the embedding layer of the LSTM using (a) 27-d physical descriptors and (b) 512-d Gaussian vectors. The number of cluster is six. (a) shows the clustered embedding of EP, which was also shown in the Figure 11.b in Kratzert et al. 2019a. (b) shows the clustering analysis of the EG-512 embedding.

streamflow. Arguably, a catchment system involves organized complexities where complexity may be manifested by randomness (Nearing et al., 2020; Dooge, 1986; Weinberg, 2001). The deep learning framework leverages this "random" complexity for streamflow
prediction and it benefits regionalization performance to a similar extent of (if not worse than) what physical descriptors provide, which is also recognized by the observation that both PEA and REA give similar streamflow prediction performance.

Although the random vector distinguishes catchments effectively and the correspond-810 ing regionalization performance is similar to physical descriptors, it is not sufficient to 811 self-prove the catchment uniqueness and it only serves as a surrogate of physical descrip-812 tors. Discussion on catchment uniqueness exceeds the scope of this paper. Our results 813 imply insights to understand hydrologic similarities within a system of catchments. We 814 analyzed and compared the patterns in the input gate (equation 9) of the EG-512 (the 815 EA-LSTM using 512-d Gaussian vector) as well as those in the EP (the EA-LSTM us-816 ing physical descriptors). This analysis intends to assess whether the original random 817 patterns of Gaussian vectors are transformed into regional patterns internally. To remind 818 the readers, the input gates are an embedding layer of 256-dimension as required by the 819 LSTM model for modeling temporal complexities, which are predetermined (See the Ap-820 pendix B of Kratzert et al. (2019a)). We conducted the K-means clustering analysis and 821 created the following map below (Figure 16). The number of clusters is set to be six as 822 suggested by Kratzert et al. (2019a). 823

It was previously shown that the EP automatically learns hydrological similarities and benefits the regionalization. The embeddings of the EP show a clear regional pattern (Figure 16 (a)). By contrast, for the Gaussian vector, the learned embeddings actually exhibit non-regional patterns, which implies the presence of hydrologic similarities that physical descriptors do not capture (Figure 16 (b)).

There are two interpretations for Figure 16. First, it clearly demonstrates that ran-829 dom vectors are not automatically transformed to physical descriptors based patterns 830 within the LSTM cell. Second, it implies that hydrologic similarities/differences are com-831 plex patterns that go beyond what physical descriptors can explain. For two catchments 832 geographically located far away from each other, certain hydrologic similarities in their 833 runoff behaviors might exist and thus cluster them into similar groups. This would not 834 be discovered by learning based on physical descriptors. Our proposed random vectors 835 seemingly allow an opportunity for deep learning models to automatically discover such 836 an implicit hydrologic similarity. The current 27-d physical descriptors might not be suf-837 ficient to fully represent cross-catchment hydrologic behavior similarities that are encoded 838 in the runoff behavior. Further research is merited to investigate catchment hydrologic 839

similarities within a catchment system, including but not limited to the role of physical descriptors in such a context.

842

## 5.4 Implications for Deep Learning in Hydrology Regionalization

Our results are delivered in a deep learning framework. The random vector approach 843 exhibits the strong modeling capacity of deep learning and shows a potential solution 844 to involving complexity into a deep learning model without explicitly incorporating hy-845 drologic processes. This approach does not add physical process understanding into the 846 model architecture; instead, it is developed purely from a data driven perspective. We 847 hypothesize that an appropriate dimension that accommodates catchment complexities 848 exists and allows deep learning models to automatically distinguish cross basin similar-849 ities and therefore benefits regional modeling. 850

Recognizing the feasibility of using random vectors when physical descriptors are 851 not available, another critical thought is that the current LSTM based models might not 852 have had an appropriate architecture to fully and explicitly leverage the physical descrip-853 tors. Physical descriptors have, by their nature, physical meanings when delivered in hy-854 drology processes. For instance, soil porosity, which is the fraction of the soil pore space, 855 impacts the amount of water stored in the vadose zone and thus it will determine the 856 amount of water released from soil into the discharge. This storage-discharge process shall 857 happen after precipitation infiltrates into the subsurface, which will often be impacted 858 by some vegetation descriptors, such as forest fraction. In short, the use of the forest frac-859 tion descriptor in process-based models comes before that of soil porosity. However, this 860 relationship is not explicitly modeled in the current LSTM-based model since both soil 861 porosity and forest fraction are used in an equal manner without any precedence to dis-862 tinguish their hydrologic roles. The same mis-utilization of other physical descriptors is 863 also present in particular between geomorphology descriptors and geology descriptors. 864 It seems natural to question the validity of how the LSTM model uses  $x^s$  since the use 865 does not explicitly reflect the hydrologic roles of physical descriptors although the LSTM 866 might implicitly learn such a relationship considering the fact that the LSTM model out-867 performs physical process based models (Kratzert et al., 2019a). Therefore, future re-868 search is merited to investigate and understand how the current LSTM uses  $x^s$ . It might 869 also improve the regionalization performance to adapt the LSTM model architecture by 870 incorporating the physical meanings of the  $x^s$  inputs. 871

The proposed random vector approach also has practical usage to assess if any given 872 physical descriptors are complete to model catchment systems complexities. Compared 873 to the performance with incomplete features (climate features, geology features, or ge-874 omorphology features) and to the performance without any physical descriptors, the pre-875 dictive performance of the random Gaussian vectors method significantly outperforms 876 in those scenarios. Random Gaussian vectors enable deep learning models to learn com-877 plexities more sufficiently than those physically limited descriptors. This insight has prac-878 tical utility for determining the sufficiency of physical descriptors in the real world, which 879 is challenging considering the uncertainties and complexities in hydrologic processes. When 880 LSTM models using a specified set of physical descriptors are outperformed by random 881 vectors, it demonstrates that those given physical descriptors are not able to resolve catch-882 ment complexities and thus suggests a need to complement them with missing features 883 for regional modeling. For instance, as a direct illustration, Figure 4 and Table 3 sug-884 gests that 27-d physical descriptors partially address hydrologic complexities and need 885 a certain degree of feature augmentation. 886

887

## 5.5 Limitations and future direction

Although the predictive performance of random vectors proves to be comparable to 27-d physical descriptors, we want to emphasize that this result is limited to gauged prediction. The deep learning model has to have training data of the basin to predict,
so the scope of this research cannot not be expanded to PUB. Therefore, recognizing this
limitation, it merits future research to leverage the complexity modeling capacity found
in random vectors into PUB. Note that the PEA (27-d physical descriptor with 512-d
embedding) is likely a suitable option to expand into PUB since it does not involve Gaussian vectors but it applicability needs further tests.

Our ability to model catchment complexities depends on the dimension of the ran-896 dom vector. We show the presence of an optimal larger d, which recognizes the existence 897 of physical processes that are not characterized, we do not provide further quantitative interpretations of the optimal d. How to utilize the observation that the optimal-d of EA-899 LSTM is 512 for regionalization models except for the embedding (section 4.6.2)? In fu-900 ture studies it will be important to identify physical processes that are not captured by 901 physical descriptors (e.g., variable recession characteristics (Beven, 2020)) and adapt ma-902 chine learning models to resolve them. Furthermore, as brought out by Beven (2020), 903 hydrological processes are unique in different catchments. Although our results do not 904 prove this statement of catchment uniqueness, the discovery of the regionalization effec-905 tiveness brought by random vectors could arguably have scientific implications in inves-906 tigating the uniqueness of catchment. Future research is merited to investigate how to 907 use random vectors to identify the uniqueness of catchments. 908

The results focus on 531 basins in the United States. Their catchment area exhibits 909 a wide range between 4 to 1980 square kilometers, which indicates strong spatial het-910 erogeneities across catchments. An interesting hypothesis to test is that a heterogeneous 911 catchment prefers high dimension Gaussian vectors to account for model complexities. 912 To test this hypothesis it will be necessary to obtain the data from catchments express-913 ing different levels of heterogeneities. Because catchments are naturally heterogeneous, 914 this test will require the use of synthetic data generated by physically based hydrolog-915 ical models. It is hypothesized that a collection of homogeneous catchment will require 916 fewer static vectors while a collection of more complex catchment will require many static 917 vectors. The synthetic data set will represent a system of catchments with a controlled 918 level of heterogeneities, which will allow an opportunity to investigate how heterogeneous 919 and homogeneous catchment systems differentiate hydrologic regionalization and mod-920 eling complexities. 921

Random vectors characterize a system of basins quite distinguishable in high di-922 mensional space. The only physically distinctive information involved becomes weather 923 inputs and associated catchment responses. This insight suggests the possibility of learn-924 ing catchment similarities from weather inputs and is thus closely related to the inverse 925 modeling problem, a field where machine learning is also advancing (Ongie et al., 2020)(Tayal 926 et al., 2022). It therefore merits future research for an improvement in unveiling catch-927 ment characterization mysteries in a physically consistent way, likely inferred from weather 928 inputs and catchment responses. 929

Our discovery has strong generalizable implications for other applications in wa-930 ter related or science problems. Regionalization can be conceptualized in a broader con-931 cept, that is, each local entity contributes to learn a regional or global model where cross 932 entity information sharing benefits the predictive performance. In the context of stream-033 flow prediction, an entity is a catchment. For water science, an entity can also be a reser-934 voir, lake, stream, etc. The target variable might vary depending on specific problems 935 to solve where each problem may require a different set of entity descriptors. Mathemat-936 ically, entities can be approximate functions in identical formulations with varying pa-937 938 rameters. The benefit of random vectors in modeling regional complexities merits further research to demonstrate their practical applicability. We expect further research can 939 test our proposed random vector approach to solve general regionalization problems across 940 disciplines. 941

## 942 6 Conclusion

953

954

In this work we showed that random vectors can be used for hydrologic regionalization when catchment physical descriptors are not available. Random vector based hydrologic regionalization shows robust performance even under data sparsity and different model strategies. This method can also identify if any given physical descriptors are sufficient to account for rainfall runoff complexities. In summary, the scientific contributions of this paper are:

- The random vector method was proposed and used for regionalization in the absence of explicit physical descriptors.
- Random vectors show robust performance even under different data sparsity scenarios and different LSTM based model selection.
  - Characterizing catchments in high-dimensional characteristics will improve regionalization performance.
- Random vectors can improve streamflow prediction when insufficient and uncertain basin characteristics are hard to distinguish basins. Thus, random vectors have a practical usage in determining if any given physical features are sufficient.

We also investigated scientific implications of the dimension of random vectors. This provides useful insights for the development of hydrologic models to address the model complexity and associated scaling issues.

## 961 Acknowledgement

The authors declare that they have no conflict of interest. This work was funded by the NSF HDR Grant: NSF Award 1934721. J.L. Nieber's effort on this project was partially supported by the USDA National Institute of Food and Agriculture, Hatch/Multistate project MN 12- 109. Access to computing facilities was provided by the Minnesota Supercomputing Institute (https://www.msi.umn.edu/). This paper relies on open-source software and all programming was done in Python environment (Van Rossum & Drake, 2009) and its relevant packages, such as pytorch (Paszke et al., 2019),

## <sup>969</sup> Data Availability Statement

Both the CAMELS data (https://doi.org/10.5065/D6G73C3Q) and the extended Maurer forcing data (https://doi.org/10.4211/hs.17c896843cf940339c3c3496d0c1c077) are publically available. The code to reproduce our work and results is available at (doi .org/10.5281/zenodo.6960022).

## 974 Author Contribution

XL and AK had the idea for Gaussian vectors. VK had the idea for one-hot vectors. All the authors were involved in the discussion of experiments design and results,
which were mainly led by XL and AK. XL conducted all the experiments and analyzed
the results together with AK. KT supervised the experiment in Section 4.6.2. XL, AK,
XJ, KC, JN, CD, MS, VK worked on the manuscripts. XL wrote the original draft and
led the editing. JN, CD supervised the manuscript from the hydrologist perspective. AK,
KC, XJ, MS, VK supervised the manuscript from the computer scientist perspective.

# <sup>982</sup> Appendix A physical descriptor description (CAMELS)

Category	Physical descrip- tors	Description
climate (9)	p_mean pet_mean aridity p_seasonality	Mean daily precipitation Mean daily potential evapotranspiration. Ratio of mean PET to mean precipitation. Seasonality and timing of precipitation. Estimated by representing annual precipitation and temperature as sine waves. Positive (negative) values indicate precip- itation peaks during the summer (winter). Values of approx. 0 indicate uniform precipitation throughout the year.
	frac_snow_daily	Fraction of precipitation falling on days with temper- atures below 0.
	high_prec_freq	Frequency of high-precipitation days ( $\geq 5$ times mean daily precipitation).
	high_prec_dur	Average duration of high-precipitation events (number of consecutive days with $\geq 5$ times mean daily precipitation).
	low_prec_freq low_prec_dur	Frequency of dry days (< 1 mm $d^{-1}$ ). Average duration of dry periods (number of consecutive days with precipitation < 1 mm $d^{-1}$ .
Geomorphology(8)	elev_mean slope_mean area_gages2 forest_frac lai_max lai_diff gvf_max gvf_diff	Catchment mean elevation. Catchment mean slope. Catchment area. Forest fraction. Maximum monthly mean of leaf area index. Difference between the max. and min. mean of the leaf area index. Maximum monthly mean of green vegetation fraction. Difference between the maximum and minimum monthly mean of the green vegetation fraction.
Geology(10)	soil_depth_pelletier soil_depth_statsgo soil_porosity soil_conductivity max_water_content sand_frac silt_frac clay_frac carb_rocks_frac	Depth to bedrock (maximum 50 m). Soil depth (maximum 1.5 m). Volumetric porosity. Saturated hydraulic conductivity. Maximum water content of the soil. Fraction of sand in the soil. Fraction of silt in the soil. Fraction of silt in the soil. Fraction of clay in the soil. Fraction of the catchment area characterized as "Carbonate sedimentary rocks".

Table A1: 27-d physical descriptors in CAMELS. Descriptions are from (Addor et al., 2017)

## 983 Appendix B FM-LSTM

FM-LSTM uses the feature modulation concept as another modelling approach. The key idea here is to use a separate gate that takes static features as input and generates a modulation vector to modulate (adapt) the features learned by a traditional LSTM. By contrast, the FM-LSTM performs the weakest modulation since the  $x^s$  is only involved for updating the hidden representation, which is the last step in an LSTM update cycle before proceeding to next timestamp.



Figure B1: FMLSTM illustration, which is based on the LSTM family illustration from "Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets" by Kratzert et al. (2019a), Hydrology and Earth System Sciences, 23, 5092 (Kratzert et al., 2019a)

As illustrated in Figure B1,  $x^s$  is mapped to an embedding layer customized for each basin (equation B6). This is then used to modulate the hidden states output (equation B7).  $x^s$  does not participate in the calculation in i[t], f[t], g[t], o[t], or c[t].

$$\boldsymbol{i}[t] = \sigma(\mathbf{W}_{\mathbf{i}}\boldsymbol{x}^{\boldsymbol{d}}[t] + \mathbf{U}_{\mathbf{i}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{i}})$$
(B1)

$$\boldsymbol{f}[t] = \sigma(\mathbf{W}_{\mathbf{f}}\boldsymbol{x}^{\boldsymbol{d}}[t] + \mathbf{U}_{\mathbf{f}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{f}})$$
(B2)

$$\boldsymbol{g}[t] = tanh(\mathbf{W}_{\mathbf{g}}\boldsymbol{x}^{\boldsymbol{d}}[t] + \mathbf{U}_{\mathbf{g}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{g}})$$
(B3)

$$\boldsymbol{o}[t] = \sigma(\mathbf{W}_{\mathbf{o}}\boldsymbol{x}^{\boldsymbol{a}}[t] + \mathbf{U}_{\mathbf{o}}\boldsymbol{h}[t-1] + \mathbf{b}_{\mathbf{o}}) \tag{B4}$$

$$\boldsymbol{c}[t] = \boldsymbol{f}[t] \odot \boldsymbol{c}[t-1] + \boldsymbol{i}[t] \odot \boldsymbol{g}[t]$$
(B5)

$$\boldsymbol{p} = \sigma(\mathbf{W}_{\mathbf{p}}\boldsymbol{x}^* + b_p) \tag{B6}$$

$$\boldsymbol{h}[t] = \boldsymbol{p} \odot \boldsymbol{o}[t] \odot tanh(\boldsymbol{c}[t]) \tag{B7}$$

Table B1: Random vector comparison in the FM-LSTM structure

models	Mean	Median
FP	0.653	0.698
FO	0.716	0.746
FG-512	0.695	0.738

For the FM-LSTM, we specify the optimal Gaussian vector dimension as the same of the EA-LSTM because they share the similar model modulation strategy, that is, static vectors enter the LSTM separately from the dynamic weather inputs. Using 27-d physical descriptors, Figure B2a illustrates that the FP yields worse prediction performance



Figure B2: Predicted performance comparison of a random vector implementation in the FM-LSTM (FO and FG-512) in contrast to the FM-LSTM using 27-d physical descriptors (FP). (a) is the NSE score scatter plot of FP and EP. (b) shows the comparison between FO and FP, while (c) shows the comparison between FG-512 and FP.

997	compared to the EP. Even so, the FM-LSTM also attains benefits performance improve
998	ment from random vectors. Both one-hot vector and Gaussian 512-d vectors lead to sig

<sup>999</sup> nificantly better predictive performance. In terms of the median, in contrast to the FP,

<sup>1000</sup> FO elevates the performance from 0.698 to 0.746, while FG-512 improves the performance

1001 to 0.738. The one-hot vector benefits are more pronounced than those of 512-d Gaus-

sian vectors in FM-LSTM. (Håkanson & Karlsson, 1984)

#### 1003 **References**

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1036

1037

1038

1039

1040

1048

- 1004Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017). The CAMELS1005data set: Catchment attributes and meteorology for large-sample studies. Hy-1006drology and Earth System Sciences, 21(10), 5293–5313. doi: 10.5194/hess-211007-5293-2017
- Alipour, M. H., & Kibler, K. M. (2018). A framework for streamflow prediction in the world's most severely data-limited regions: Test of applicability and performance in a poorly-gauged region of China. Journal of Hydrology, 557, 41–54.
   Retrieved from https://doi.org/10.1016/j.jhydrol.2017.12.019 doi: 10.1016/j.jhydrol.2017.12.019
- Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Miralles, D. G., McVicar, T. R.,
  Schellekens, J., & Bruijnzeel, L. A. (2016). Global-scale regionalization of
  hydrologic model parameters. Water Resources Research, 52(5), 3599-3622.
  Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/
  10.1002/2015WR018247 doi: https://doi.org/10.1002/2015WR018247
  - Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2), 157-166. doi: 10.1109/72.279181
  - Besaw, L. E., Rizzo, D. M., Bierman, P. R., & Hackett, W. R. (2010). Advances in ungauged streamflow prediction using artificial neural networks. *Journal of Hydrology*, 386(1-4), 27–37. Retrieved from http://dx.doi.org/10.1016/j
    .jhydrol.2010.02.037 doi: 10.1016/j.jhydrol.2010.02.037
  - Beven, K. (1987). Towards a new paradigm in hydrology. Water for the future. Proc. Rome symposium, 1987(164), 393–403.
  - Beven, K. (1989). CHANGING IDEAS IN HYDROLOGY- THE CASE OF PHYSICALLY-BASED MODELS., 105, 157–172.
  - Beven, K. (2001). How far can we go in distributed hydrological modelling? *Hydrology and Earth System Sciences*, 5(1), 1–12. doi: 10.5194/hess-5-1-2001
- Beven, K. (2002). Towards an alternative blueprint for a physically based digitally
   simulated hydrologic response modelling system. *Hydrological Processes*, 16(2),
   189–206. doi: 10.1002/hyp.343
- Beven, K. (2020). Deep learning, hydrological processes and the uniqueness of place.
   Hydrological Processes, 34 (16), 3608–3613. doi: 10.1002/hyp.13805
  - Beven, K., & Binley, A. (1992). The future of distributed models: Model calibration and uncertainty prediction. *Hydrological Processes*, 6(3), 279–298. Retrieved from http://dx.doi.org/10.1002/hyp.3360060305
  - Blöschl, G., & Sivapalan, M. (1995). Scale issues in hydrological modelling: A review. Hydrological Processes, 9(3-4), 251–290. doi: 10.1002/hyp.3360090305
- Burnash, R. J. C. (1995). The NWS river forecast system-catchment modeling.
   *Computer models of watershed hydrology*, 311–366.
- Choubin, B., Solaimani, K., Rezanezhad, F., Habibnejad Roshan, M., Malekian,
   A., & Shamshirband, S. (2019). Streamflow regionalization using a similarity approach in ungauged basins: Application of the geo-environmental
   signatures in the Karkheh River Basin, Iran. Catena, 182(June), 104128.
   Retrieved from https://doi.org/10.1016/j.catena.2019.104128 doi:
- de Lavenne, A., Andréassian, V., Thirel, G., Ramos, M. H., & Perrin, C. (2019). A
   Regularization Approach to Improve the Sequential Calibration of a Semidis tributed Hydrological Model. Water Resources Research, 55 (11), 8821–8839.
   doi: 10.1029/2018WR024266

10.1016/j.catena.2019.104128

Dooge, J. C. I. (1986). Looking for hydrologic laws. Water Resources Research,
 22(9S), 46S-58S. Retrieved from https://agupubs.onlinelibrary.wiley

1055	.com/doi/abs/10.1029/WR022i09Sp0046S doi: https://doi.org/10.1029/
1056	WR022i09Sp0046S
1057	Drost, N. F. WA. A., S., & Mudersbach, C. (2021). The impact of land cover data
1058	on rainfall-runoff prediction using an entity-aware-lstm doi: https://doi.org/
1059	10.5194/egusphere-egu21-1136,2021.
1060	Ecrepont, S., Cudennec, C., Anctil, F., & Jaffrézic, A. (2019). PUB in Québec: A
1061	robust geomorphology-based deconvolution-reconvolution framework for the
1062	spatial transposition of hydrographs. <i>Journal of Hydrology</i> , 570(January),
1063	378–392. doi: 10.1016/j.jhydrol.2018.12.052
1064	Feng, D., Fang, K., & Shen, C. (2020). Enhancing Streamflow Forecast and Ex-
1065	tracting Insights Using Long-Short Term Memory Networks With Data Inte-
1066	gration at Continental Scales. Water Resources Research, 56(9), 1–24. doi:
1067	10.1029/2019WR026793
1068	Frame, J., Nearing, G., Kratzert, F., Raney, A., & Rahman, M. (2020, Jul).
1069	Post processing the u.s. national water model with a long short-term mem-
1070	ory network. EarthArXiv. Retrieved from eartharxiv.org/4xhac doi:
1071	10.31223/osf.io/4xhac
1072	Freeze, R. A. (1974). Streamflow generation. <i>Reviews of Geophysics</i> , 12(4), 627–647.
1073	doi: 10.1029/RG012i004p00627
1074	Freeze, R. A., & Harlan, R. (1969). BLUEPRINT FOR A PHYSICALLY-BASED.
1075	DIGITALLY-SIMULATED HYDROLOGIC RESPONSE MODEL. Journal of
1076	Hydrology, 9, 237–258.
1077	Gauch, M., Mai, J., & Lin, J. (2021). The proper care and feeding of CAMELS:
1078	How limited training data affects streamflow prediction. <i>Environmental Mod</i> -
1079	elling and Software, 135, 0–2, doi: 10.1016/i.envsoft.2020.104926
1080	Gauch, M., et al. (2021). The proper care and feeding of camels: How limited
1081	training data affects streamflow prediction. Environmental Modelling Soft-
1082	ware, 135, 104926. Retrieved from https://www.sciencedirect.com/
1083	science/article/pii/S136481522030983X doi: https://doi.org/10.1016/
1084	i.envsoft.2020.104926
1085	Guo, Y., Zhang, Y., Zhang, L., & Wang, Z. (2021). Regionalization of hydrolog-
1086	ical modeling for predicting streamflow in ungauged catchments: A compre-
1087	hensive review. Wiley Interdisciplinary Reviews: Water, 8(1), 1–32. doi:
1088	10.1002/wat2.1487
1089	Hochreiter, S., & Schmidhuber, J. (1997, 11), Long Short-Term Memory, Neu-
1090	ral Computation, 9(8), 1735-1780. Retrieved from https://doi.org/10.1162/
1091	neco.1997.9.8.1735 doi: 10.1162/neco.1997.9.8.1735
1092	Hsu, K. Gupta, H. V., & Sorooshian, S. (1995). Artificial Neural Network Modeling
1093	of the Rainfall-Runoff Process. Water Resources Research, 31(10), 2517–2530.
1094	doi: 10.1029/95WR01955
1095	Håkanson, L., & Karlsson, B. (1984). On the relationship between regional ge-
1096	omorphology and lake morphometry—a swedish example. Geografiska An-
1097	naler: Series A. Physical Geography, 66(1-2), 103-119. Retrieved from
1098	https://doi.org/10.1080/04353676.1984.11880102 doi: 10.1080/
1099	04353676.1984.11880102
1100	Kingma, D. P., & Ba, J. (2017). Adam: A method for stochastic ontimization.
1101	Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall –
1102	runoff modelling using Long Short-Term Memory (LSTM) networks. Hudrol-
1103	oay and Earth System Sciences, 22, 6005–6022.
1104	Kratzert, F., Klotz, D., Herrnegger, M. Sampson, A.K. Hochreiter, S. & Nearing
1105	G. S. (2019b). Toward improved predictions in ungauged basins: Exploiting
1106	the power of machine learning. Water Resources Research, 55(12) 11344-
1107	11354. Retrieved from https://agupubs.onlinelibrary.wilev.com/doi/
1108	abs/10.1029/2019WR026065 doi: https://doi.org/10.1029/2019WR026065
1109	Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G.
	, , , , , , , , , , , , , , , , , , ,

1110	(2019a). Towards learning universal, regional, and local hydrological behaviors
1111	via machine learning applied to large-sample datasets. Hydrology and Earth System Sciences $23(12)$ 5089–5110 doi: 10.5194/bess-23-5089-2019
1112	Ling Y Lattermain D P Wood F F & Burges S L $(1004)$ A simple by
1113	drologically based model of land surface water and energy flyes for general
1114	airculation models $I_{aurral}$ of Coonbusical Research: Atmospheres $00(D7)$
1115	14415 14428 Detrieved from https://orupuba.onlinelibrory.uiley.com/
1116	doi /obg/10_1020/04 ID00482_doi: https://doi.org/10_1020/04 ID00483
1117	$M_{0} K = E_{0} m_{0} D = L_{0} m_{0} K = M_{0} D = M_{0} M_{0} K = M_{0} D = M_{0} M_{0} K = M_{0} D = M_{0} M_$
1118	Ma, K., Feng, D., Lawson, K., Isal, W., Llang, C., Huang, A., Shen, C. (2021).
1119	to improve hydrologic usits prediction in data sparse regions
1120	Research 1–26 doi: 10.1020/2020wr028600
1121	Mauror F P Wood A W Adam I C Lattenmaior D P & Nijsson B
1122	(2002) A long-term hydrologically based dataset of land surface fluxes
1123	and states for the conterminous united states
1124	3237 - 3251 Betrieved from https://iournals.ametsoc.org/vieu/
1125	$\frac{1}{10000000000000000000000000000000000$
1120	doi: 10.1175/1520-0442/2002)015/3237·ALTHBD\2.0.CO-2
1127	McDonnell I I Sivapalan M Vaché K Dunn S Grant G Haggerty B
1128	Weiler M (2007) Moving beyond heterogeneity and process complexity: A
1129	new vision for watershed hydrology Water Resources Research 13(7) 1-6
1121	doi: 10.1029/2006WB005467
1122	Nearing C. S. Kratzert F. Sampson A. K. Craig S. Frame I. M. Klotz D.
1132	& Gupta H V (2020) What Bole Does Hydrological Science Play in
1133	the Age of Machine Learning? Water Resources Research 1–17 doi:
1135	10.31223/osf.io/3sx6g
1126	Ongie G. Jalal A. Metzler, C. A. Baraniuk, B. G. Dimakis, A. G. & Willett, B.
1130	(2020) Deen learning techniques for inverse problems in imaging
1157	(2020). Deep tearning teening all for interior proteine in intaging.
1138	Pagliero L. Bouraoui F. Diels J. Willems P. & McIntyre N. (2019) Investigat-
1138	Pagliero, L., Bouraoui, F., Diels, J., Willems, P., & McIntyre, N. (2019). Investigat- ing regionalization techniques for large-scale hydrological modelling. <i>Journal of</i>
1138 1139 1140	Pagliero, L., Bouraoui, F., Diels, J., Willems, P., & McIntyre, N. (2019). Investigat- ing regionalization techniques for large-scale hydrological modelling. <i>Journal of</i> <i>Hudrologu</i> , 570 (September 2017), 220–235. Retrieved from https://doi.org/
1138 1139 1140 1141	Pagliero, L., Bouraoui, F., Diels, J., Willems, P., & McIntyre, N. (2019). Investigat- ing regionalization techniques for large-scale hydrological modelling. <i>Journal of</i> <i>Hydrology</i> , 570 (September 2017), 220–235. Retrieved from https://doi.org/ 10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071
1138 1139 1140 1141	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. <i>Journal of Hydrology</i>, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., M., Chin-</li> </ul>
1138 1139 1140 1141 1142 1143	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep</li> </ul>
1138 1139 1140 1141 1142 1143 1144	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Bevgelzimer, F. d'Alché-</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information process-</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Re-</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research,</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N.,</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1145 1147 1148 1149 1150 1151 1152 1153	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales.</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323-4346. doi:</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220-235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024-8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323-4346. doi: 10.5194/hess-2017-89</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364–4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323–4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lak-</li> </ul>
1138 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364–4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323–4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Un-</li> </ul>
1138 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323-4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Ungauged Basins (PUB), 2003-2012: Shaping an exciting future for the hy-</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364–4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323–4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Ungauged Basins (PUB), 2003-2012: Shaping an exciting future for the hydrological sciences. Hydrological Sciences Journal, 48(6), 857–880. doi:</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1155 1156 1157 1158 1159 1160 1161	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55 (5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323-4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Ungauged Basins (PUB), 2003-2012: Shaping an exciting future for the hydrological sciences. Hydrological Sciences Journal, 48(6), 857-880. doi: 10.1623/hysj.48.6.857.51421</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1155 1156 1157 1158 1159 1160 1161 1161	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364-4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323-4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Ungauged Basins (PUB), 2003-2012: Shaping an exciting future for the hydrological sciences. Hydrological Sciences Journal, 48(6), 857-880. doi: 10.1623/hysj.48.6.857.51421</li> <li>Tayal, K., Jia, X., Ghosh, R., Willard, J., Read, J., &amp; Kumar, V. (2022). Invert-</li> </ul>
1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1161 1162 1162 1163	<ul> <li>Pagliero, L., Bouraoui, F., Diels, J., Willems, P., &amp; McIntyre, N. (2019). Investigating regionalization techniques for large-scale hydrological modelling. Journal of Hydrology, 570 (September 2017), 220–235. Retrieved from https://doi.org/10.1016/j.jhydrol.2018.12.071 doi: 10.1016/j.jhydrol.2018.12.071</li> <li>Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Chintala, S. (2019). Pytorch: An imperative style high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, &amp; R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Curran Associates, Inc. Retrieved from http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style -high-performance-deep-learning-library.pdf</li> <li>Prieto, C., Le Vine, N., Kavetski, D., García, E., &amp; Medina, R. (2019). Flow Prediction in Ungauged Catchments Using Probabilistic Random Forests Regionalization and New Statistical Adequacy Tests. Water Resources Research, 55(5), 4364–4392. doi: 10.1029/2018WR023254</li> <li>Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Attinger, S. (2017). Toward seamless hydrologic predictions across scales. Hydrology and earth system sciences discussions, 2017(89), 4323–4346. doi: 10.5194/hess-2017-89</li> <li>Sivapalan, M., Takeuchi, K., Franks, S. W., Gupta, V. K., Karambiri, H., Lakshmi, V., Zehe, E. (2003). IAHS Decade on Predictions in Ungauged Basins (PUB), 2003-2012: Shaping an exciting future for the hydrological sciences. Hydrological Sciences Journal, 48(6), 857–880. doi: 10.1623/hysj.48.6.857.51421</li> <li>Tayal, K., Jia, X., Ghosh, R., Willard, J., Read, J., &amp; Kumar, V. (2022). Invertibility aware integration of static and time-series data: An application to lake</li> </ul>

1165	ence on data mining (sdm).
1166	Thornton, M., Shrestha, R., Wei, Y., Thornton, P., Kao, S., & Wilson, B. (2020).
1167	Daymet: Daily surface weather data on a 1-km grid for north america,
1168	version 4. ORNL Distributed Active Archive Center. Retrieved from
1169	https://daac.ornl.gov/cgi-bin/dsviewer.pl?ds_id=1840 doi:
1170	10.3334/ORNLDAAC/1840
1171	Van Rossum, G., & Drake, F. L. (2009). Python 3 reference manual. Scotts Valley,
1172	CA: CreateSpace.
1173	Weinberg, G. M. (2001). An introduction to general systems thinking (silver anniver-
1174	sary ed.). USA: Dorset House Publishing Co., Inc.
1175	Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Mocko,
1176	D. (2012). Continental-scale water and energy flux analysis and valida-
1177	tion for the north american land data assimilation system project phase $2$
1178	(nldas-2): 1. intercomparison and application of model products. Journal
1179	of Geophysical Research: Atmospheres, 117(D3). Retrieved from https://
1180	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2011JD016048 doi:
1181	https://doi.org/10.1029/2011JD016048
1182	Zamoum, S., & Souag-Gamane, D. (2019). Monthly streamflow estimation in
1183	ungauged catchments of northern Algeria using regionalization of concep-
1184	tual model parameters. Arabian Journal of Geosciences, $12(11)$ . doi:
1185	10.1007/s12517-019-4487-9

-37-

# Supporting Information for "Regionalization in a global hydrologic deep learning model: from physical descriptors to random vectors"

Xiang Li<sup>1</sup>, Ankush Khandelwal<sup>2</sup>, Xiaowei Jia<sup>3</sup>, Kelly Cutler<sup>2</sup>, Rahul Ghosh<sup>2</sup>,

Arvind Renganathan<sup>2</sup>, Shaoming Xu<sup>2</sup>, Kshitij Tayal<sup>2</sup>, John Nieber<sup>1</sup>,

Christopher Duffy<sup>4</sup>, Michael Steinbach<sup>2</sup>, Vipin Kumar<sup>2</sup>

<sup>1</sup>Department of Bioproducts and Biosystems Engineering, University of Minnesota Twin Cities, St.Paul, MN, USA <sup>2</sup>Department of Computer Science and Engineering, University of Minnesota Twin Cities, Minneapolis, MN, USA

 $^3{\rm School}$  of Computing and Information, University of Pittsburgh, Pittsburgh, PA, USA

<sup>4</sup>Department of Civil and Environmental Engineering, Pennsylvania State University, State College, PA, USA

## Contents of this file

1. Tables S1 to S7  $\,$ 

## Introduction

This supplement information document report the p-value statistics regarding all NSE comparison in the paper — "Regionalization in a global hydrologic deep learning model: from physical descriptors to random vectors". The captions for each p-value table below marks the corresponding NSE comparison table in the paper (from Table 3 to Table 10). All "NA"s mark the corresponding the p-value entries comparing with itself and thus not able to calculate it and reports "NA".

July 31, 2022, 11:04pm

# Tables S1 to S7

Table S1. The p-value table for Table 3 and Table 10 in the main article.

Catchment static vectors	p-value (to EP)	p-value (to EG-512)
27-d physical vectors (EP)	NA	9.25e-06
512-d random vectors (EG-512)	9.25e-06	NA
one-hot vectors (EO)	3.14e-10	0.121
27-d physical descriptors + 512-d embedding (PEA)	1.40e-10	0.003
27-d Gaussian vectors + 512-d embedding (REA)	1.41e-09	2.09e-05

Table S2. The p-value table for Table 4 in the main article

k-basin group		10	50	100
	d	p-value to	o its EP co	ounterpart
Gaussian vector (EG-d)	2	1.41e-56	4.05e-58	1.99e-59
	8	5.29e-05	6.69e-51	2.10e-43
	16	2.29e-19	7.47e-35	1.72e-21
	32	2.17e-07	6.31e-06	0.003
	64	1.49e-09	1.08e-17	8.84e-27
	128	7.95e-55	2.66e-58	6.13e-42
	256	2.17e-72	9.55e-65	1.29e-53
	512	4.18e-63	8.24e-59	1.64e-52
	1024	9.46e-31	2.28e-33	2.98e-36
one-hot (EO)		6.05e-45	8.24e-40	1.47e-24
27-d physical descriptors	(EP)	NA	NA	NA

Table S3.	The p-value	table for	Table 5	5 in	the	$\operatorname{main}$	article
-----------	-------------	-----------	---------	------	-----	-----------------------	---------

1				
Number of training ye	1	2	5	
$\overline{\text{Gaussian 512-d (EG-512)}}$	to EP	2.69e-09	0.022	1.67e-07
one-hot $(EO)$	to EP	0.049	1.65e-36	1.42e-18

Table S4. The p-value table for Table 6 in the main article.

models	p-value (to CP)
CP	NA
CO	9.72e-05
CG-16	0.003
CG-32	2.37e-06

1				
k-basin group		10	50	100
	d	p-value to	its CP co	unterpart
Gaussian vector (CG-d)	2	2.25e-48	1.20e-56	1.05e-54
	8	4.06e-48	3.34e-30	1.02e-24
	16	6.64 e- 36	2.13e-14	2.29e-08
	32	1.07e-16	1.01e-05	0.231
	64	7.970e-18	0.0112	0.956
	128	9.72e-40	1.89e-14	0.027
	256	1.78e-64	2.57e-51	1.11e-29
	512	9.23e-73	6.91e-79	7.02e-72
	1024	1.68e-78	2.68e-83	5.50e-82
one-hot (CO)		3.33e-32	3.93e-05	0.173
27-d physical descriptors	(CP)	NA	NA	NA

Table S5. The p-value table for Table 7 in the main article

Table S6. The p-value table for Table 8 in the main article

catchment static vectors	p-value (to EP)
(0-d) Without static features	1.09e-75
512-d Gaussian vectors (EG- $512$ )	9.25e-06
9-d climate features	4.21e-72
10-d geology features	7.04e-62
8-d geomorphology features	1.87e-62
27-d physical descriptors (EP)	NA

Table S7. The p-value table for Table 9 in the main article

catchment static vectors	d	p-value (to EP)	p-value (to EG-d, the same d)
mixed Gaussian vector (EM-d)	32	1.55e-08	6.38e-03
	64	2.42e-04	4.75e-17
	128	2.21e-08	3.47e-06
	256	2.44e-08	2.43e-04
	512	5.63e-14	2.81e-07
	1024	3.31e-04	7.19e-11