



# Water Resources Research

## TECHNICAL REPORTS: METHODS

10.1029/2019WR026065

### Special Section:

Big Data & Machine Learning in Water Sciences: Recent Progress and Their Use in Advancing Science

### Key Points:

- Overall accuracy of LSTMs in ungauged basins is comparable to standard hydrology models in gauged basins
- There is sufficient information in catchment characteristics data to differentiate between catchment-specific rainfall-runoff behaviors

### Correspondence to:

G. S. Nearing,  
gsnearing@ua.edu

### Citation:

Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019). Toward improved predictions in ungauged basins: Exploiting the power of machine learning. *Water Resources Research*, 55, 11, 344–11,354.  
<https://doi.org/10.1029/2019WR026065>

Received 29 JUL 2019

Accepted 19 NOV 2019

Accepted article online 23 NOV 2019

Published online 23 DEC 2019

## Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine Learning

Frederik Kratzert<sup>1</sup>, Daniel Klotz<sup>1</sup>, Mathew Herrnegger<sup>2</sup>, Alden K. Sampson<sup>3</sup>, Sepp Hochreiter<sup>1</sup>, and Grey S. Nearing<sup>4</sup>

<sup>1</sup>LIT AI Lab and Institute for Machine Learning, Johannes Kepler University, Linz, Austria, <sup>2</sup>Institute for Hydrology and Water Management, University of Natural Resources and Life Sciences, Vienna, Austria, <sup>3</sup>Upstream Tech, Natel Energy Inc., Alameda, CA, USA, <sup>4</sup>Department of Geological Sciences, University of Alabama, Tuscaloosa, AL, USA

**Abstract** Long short-term memory (LSTM) networks offer unprecedented accuracy for prediction in ungauged basins. We trained and tested several LSTMs on 531 basins from the CAMELS data set using k-fold validation, so that predictions were made in basins that supplied no training data. The training and test data set included ~30 years of daily rainfall-runoff data from catchments in the United States ranging in size from 4 to 2,000 km<sup>2</sup> with aridity index from 0.22 to 5.20, and including 12 of the 13 IGPB vegetated land cover classifications. This effectively “ungauged” model was benchmarked over a 15-year validation period against the Sacramento Soil Moisture Accounting (SAC-SMA) model and also against the NOAA National Water Model reanalysis. SAC-SMA was calibrated separately for each basin using 15 years of daily data. The out-of-sample LSTM had higher median Nash-Sutcliffe Efficiencies across the 531 basins (0.69) than either the calibrated SAC-SMA (0.64) or the National Water Model (0.58). This indicates that there is (typically) sufficient information in available catchment attributes data about similarities and differences between catchment-level rainfall-runoff behaviors to provide out-of-sample simulations that are generally more accurate than current models under ideal (i.e., calibrated) conditions. We found evidence that adding physical constraints to the LSTM models might improve simulations, which we suggest motivates future research related to physics-guided machine learning.

### 1. Introduction

Science and society are firmly in the age of machine learning (ML; McAfee & Brynjolfsson, 2017). ML models currently outperform state-of-the-art techniques at some of the most sophisticated domain problems across the Natural Sciences (e.g., AlQuraishi, 2019; He et al., 2019; Liu et al., 2016; Mayr et al., 2016). In Hydrology, the first demonstration of ML outperforming a process-based model that we are aware of was by Hsu et al. (1995), who compared a calibrated Sacramento Soil Moisture Accounting Model (SAC-SMA) against a feed-forward artificial neural network across a range of flow regimes. More recently, Nearing et al. (2018) compared neural networks against the half-hourly surface energy balance of hydrometeorological models used operationally by several international weather and climate forecasting agencies, and showed that the former generally out-performed the latter at out-of-sample FluxNet sites. In a companion paper to this one, Kratzert et al. (2019) showed that regionally trained long short-term memory (LSTM) network outperforms basin-specific calibrations of several traditional hydrology models and demonstrated that LSTM-type models were able to extract information from observable catchment characteristics to differentiate between different rainfall-runoff behaviors in hydrologically diverse catchments. The purpose of this paper is to show that we can leverage this capability for prediction in ungauged basins.

There is a long-standing discussion in the field of Hydrology about the relative merits of data-driven versus process-driven models (e.g., Klemeš, 1986). In their summary of a recent workshop on “Big Data and the Earth Sciences,” Sellars (2018) noted that “Many participants who have worked in modeling physical-based systems continue to raise caution about the lack of physical understanding of machine learning methods that rely on data-driven approaches.” It is often argued that data-driven models might underperform relative to models that include explicit process representations in conditions that are dissimilar to training data (e.g., Kirchner, 2006; Milly et al., 2008; Vaze et al., 2015). While this may or may not be true (we are unaware of any study that has tested this hypothesis directly), in any case where an ML model *does* outperform relative to a given process-based model, we can conclude that the process-based model does not take advantage of

©2019. The Authors.

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

the full information content of the input/output data (Nearing & Gupta, 2015). At the very least, such cases indicate that there is potential to improve the process-based model(s).

One of the situations where the accuracy of out-of-sample predictions matter is for prediction in ungauged basins (PUB). PUB was the decadal problem of the International Association of Hydrological Sciences (IAHS) from 2003–2012 (Hrachowitz et al., 2013; Sivapalan et al., 2003). State-of-the-art regionalization, parameter transfer, catchment similarity, and surrogate basin techniques (e.g., Parajka et al., 2013; Razavi & Coulibaly, 2012; Samaniego et al., 2017) result in streamflow predictions that are less accurate than from models calibrated individually in gauged catchments. Current community best practices for PUB center around obtaining detailed local knowledge of a particular basin (Blöschl, 2016), which is expensive for individual catchments and impossible for large-scale (e.g., continental) simulations like those from the U.S. National Water Model (NWM; Salas et al., 2018) or the streamflow component of the North American Land Data Assimilation System (NLDAS; Xia et al., 2012). Moreover, Vrugt et al. (2006) argued that reliable streamflow predictions from lumped catchment models typically require at least 2 to 3 years of gauge data for calibration (even this is likely an underestimate of the amount of data necessary for reliable model calibration). PUB remains an important challenge because the majority of streams in the world are either ungauged or poorly gauged (Goswami et al., 2007; Sivapalan, 2003), and the number of gauged catchments, even in the United States, is shrinking (Fekete et al., 2015).

In this technical note, we demonstrate an ML strategy for PUB. Our results show that out-of-sample LSTMs outperform, on average, a conceptual model (SAC-SMA) calibrated independently for each catchment, and also a distributed, process-based model (NWM). The purpose of this demonstration is twofold. First, to show that there is sufficient information in the available hydrological data record to provide meaningful predictions in ungauged basins—at least a significant portion of the time. Second, to show that ML offers a promising path forward for extracting this information, and for PUB in general. The current authors are unaware of any existing model that performs as well, on average, as the LSTMs that we demonstrate here. At the end of this technical note we offer some thoughts—both philosophical and practical—about future work that could be done to advance the utility of ML in a complex systems science like Hydrology.

To reemphasize our primary findings succinctly, ML in *ungaaged basins* outperforms, on average (i.e., in more catchments than not) a lumped conceptual model calibrated in *gauged basins*, and also a state-of-the-art distributed process-based model. This rapid correspondence is intended to highlight initial results that might motivate continued development of these and similar techniques—this is not intended to be a comprehensive analysis of the application of LSTMs or deep learning in general to PUB.

## 2. Data

Experimental data for our analysis came from the publicly available Catchment Attributes and Meteorology for Large-Sample Studies (CAMELS) data set curated by National Center for Atmospheric Research (NCAR; Addor et al., 2017; Newman et al., 2014; Newman et al., 2015). CAMELS consists of 671 catchments in the continental United States ranging in size from 4 to 25,000 km<sup>2</sup>. These catchments were chosen from the available gauged catchments in the United States due to the fact that they are largely natural and have long gauge records (1980–2010) available from the United States Geological Survey National Water Information System. CAMELS includes daily forcing from Daymet, Maurer, and NLDAS, as well as several static catchment attributes related to soils, climate, vegetation, topography, and geology (Addor et al., 2018). It is important to point out that these catchment attributes were derived from maps, remote sensing products, and climate data that are generally available over the continental United States and either exactly or in close approximation, globally. For this project, we used only 531 of 671 CAMELS catchments—these were the same basins that were used for model benchmarking by Newman et al. (2017), who removed basins from the full CAMELS data set with (i) large discrepancies between different methods of calculating catchment area and (ii) areas larger than 2,000 km<sup>2</sup>.

The CAMELS repository also includes daily streamflow values simulated by 10 SAC-SMA models calibrated separately in each catchment using Shuffled Complex Evolution (SCE; Duan et al., 1993) with 10 random seeds. Each SAC-SMA was calibrated on 15 years of data in each catchment (1980–1995). These calibrations were performed in previous work by NCAR (Newman et al., 2015). We used this ensemble of SAC-SMA models as a benchmark for our LSTMs. In addition, we benchmarked against the NWM reanalysis, which spans the years 1993–2017 (<https://docs.opendata.aws/nwm-archive>). All performance statistics that we report

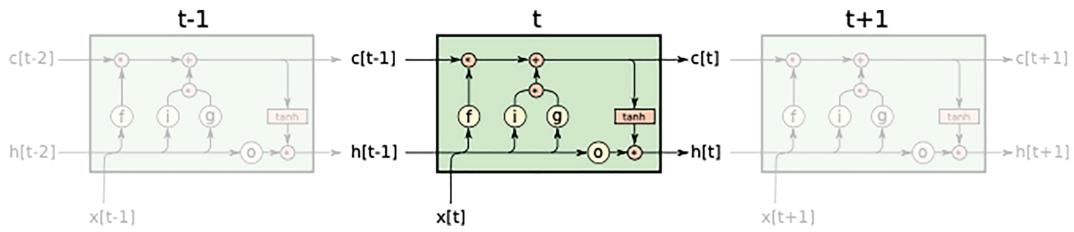


Figure 1. Visualization of (a) the standard LSTM cell as defined by equations (1)–(6).

(for all models) are from the water years 1996–2010, so that the SAC-SMA models were tested out of sample in time but at the same basins where they were calibrated.

### 3. Methods

#### 3.1. A Brief Overview of LSTM Networks

LSTMs are a type of recurrent neural network (RNN) first proposed by Hochreiter and Schmidhuber (1997). LSTMs have memory cells that are analogous to the states of a traditional dynamical systems model, which make them useful for simulating natural systems like watersheds. Compared with other types of recurrent neural networks, LSTMs avoid exploding and/or vanishing gradients, which allows them to learn long-term dependencies between input and output features. This is desirable for modeling catchment processes like snow accumulation and seasonal vegetation patterns that have relatively long timescales as compared with input-driven processes like direct surface runoff. Kratzert, Klotz, et al. (2018) applied LSTMs to the problem of rainfall-runoff modeling and later demonstrated that the internal memory states of the network were highly correlated with observed snow and soil moisture states without the model seeing any type of snow or soil moisture data during training (Kratzert, Herrnegger, Kratzert et al., 2018).

Figure 1 provides an illustration of an LSTM, which works as follows. The model takes a time series (more generally, a sequence) of inputs  $\mathbf{x} = [\mathbf{x}[1], \dots, \mathbf{x}[T]]$  of data over  $T$  time steps, where each element  $\mathbf{x}[t]$  is a vector containing features (model inputs) at time step  $t$ . This is not dissimilar to any standard hydrological simulation model (i.e., is it not a one-step-ahead forecast model). The LSTM model structure is described by the following equations:

$$\mathbf{x}[t] + \mathbf{U}_i \mathbf{h}[t-1] + \mathbf{b}_i \quad (1)$$

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

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

$$\mathbf{o}[t] = \sigma(\mathbf{W}_o \mathbf{x}[t] + \mathbf{U}_o + \mathbf{b}_o) \quad (4)$$

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

$$\mathbf{h}[t] = \mathbf{o}[t] \odot \tanh(\mathbf{c}[t]), \quad (6)$$

where  $\mathbf{i}[t]$ ,  $\mathbf{f}[t]$ , and  $\mathbf{o}[t]$  are the *input gate*, *forget gate*, and *output gate*, respectively,  $\mathbf{g}[t]$  is the *cell input* and  $\mathbf{x}[t]$  is the *network input* at time step  $t$  ( $1 \leq t \leq T$ ),  $\mathbf{h}[t-1]$  is the *recurrent input*  $\mathbf{c}[t-1]$  the *cell state* from the previous time step. At the first time step, the hidden and cell states are initialized as a vector of zeros.  $\mathbf{W}$ ,  $\mathbf{U}$ , and  $\mathbf{b}$  are calibrated parameters. These are specific to each gate, and subscripts indicate which gate the particular weight matrix/vector is associated with.  $\sigma(\cdot)$  is the sigmoid activation function,  $\tanh(\cdot)$  the hyperbolic tangent function, and  $\odot$  is element-wise multiplication. The intuition is that the cell states ( $\mathbf{c}[t]$ ) characterize the memory of the system. These are modified by (i) the forget gate ( $\mathbf{f}[t]$ ), which allows attenuation of information in the states over time, and by (ii) a combination of the input gate ( $\mathbf{i}[t]$ ) and cell update ( $\mathbf{g}[t]$ ), which can add new information. In the latter case, the cell update contains information to be added to each cell state, and the input gate (which is a sigmoid function) controls which cells are “allowed”

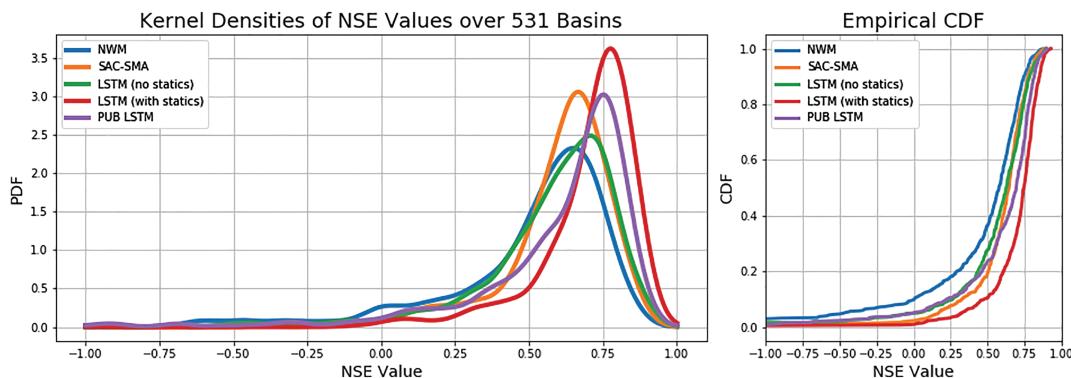
**Table 1**  
*Table of LSTM Inputs*

Meteorological forcing data	
Maximum air temp	2 m daily maximum air temperature (°C)
Minimum air temp	2 m daily minimum air temperature (°C)
Precipitation	Average daily precipitation (mm/day)
Radiation	Surface-incident solar radiation (W/m <sup>2</sup> )
Vapor pressure	Near-surface daily average (P <sub>a</sub> )
Static catchment attributes	
Precipitation mean	Mean daily precipitation.
PET mean	Mean daily potential evapotranspiration
Aridity index	Ratio of Mean PET to Mean Precipitation
Precip seasonality	Estimated by representing annual precipitation and temperature as sin waves Positive (negative) values indicate precipitation peaks during the summer (winter). Values of ~0 indicate uniform precipitation throughout the year.
Snow fraction	Fraction of precipitation falling on days with temp < 0 °C.
High precipitation frequency	Frequency of days with ≤ 5× mean daily precipitation
High precip duration	Average duration of high precipitation events (number of consecutive days with ≤ 5× mean daily precipitation).
Low precip frequency	Frequency of dry days (< 1 mm/day).
Low precip duration	Average duration of dry periods (number of consecutive days with precipitation < 1 mm/day).
Elevation	Catchment mean elevation.
Slope	Catchment mean slope.
Area	Catchment area.
Forest fraction	Fraction of catchment covered by forest.
LAI max	Maximum monthly mean of leaf area index.
LAI difference	Difference between the max. and min. mean of the leaf area index.
GVF max	Maximum monthly mean of green vegetation fraction.
GVF difference	Difference between the maximum and minimum monthly mean of the green vegetation fraction.
Soil depth (Pelletier)	Depth to bedrock (maximum 50 m).
Soil depth (STATSGO)	Soil depth (maximum 1.5 m).
Soil Porosity	Volumetric porosity.
Soil conductivity	Saturated hydraulic conductivity.
Max water content	Maximum water content of the soil.
Sand fraction	Fraction of sand in the soil.
Silt fraction	Fraction of silt in the soil.
Clay fraction	Fraction of clay in the soil.
Carbonate rocks fraction	Fraction of the catchment area characterized as “carbonate sedimentary rocks.”
Geological permeability	Surface permeability (log10).

to receive new information. Finally, the output gate ( $\mathbf{o}[t]$ ) controls the flow of information from states to model output.

### 3.2. Experimental Design

The LSTMs used in this study took as inputs at each time step the NLDAS meteorological forcing data listed in Table 1. Additionally, at each time step, the meteorological inputs were augmented with the catchment



**Figure 2.** Frequencies of NSE values from 531 catchments given by “gauged” and “ungauged” LSTMs, calibrated (gauged) SAC-SMA, and the National Water Model reanalysis.

attributes also listed in Table 1. These catchment attributes were described in detail by Addor et al. (2017) and remain constant in time throughout the simulation (training and testing). In total we used 32 LSTM inputs at each daily time step: 5 meteorological forcings and 27 catchment characteristics. All LSTMs were configured to have 256 cell states with a dropout rate of 0.4 applied to the LSTM output before a single regression layer.

We trained and tested three types of LSTM models:

1. *Global LSTM without static features*: LSTMs with only meteorological forcing inputs, and without catchment attributes, trained on all catchments simultaneously (without k-fold validation).
2. *Global LSTM with static features*: LSTMs with both meteorological forcings and catchment characteristics as inputs, trained on all catchments simultaneously (without k-fold validation).
3. *PUB LSTM*: LSTMs with both meteorological forcings and catchment characteristics as inputs, trained and tested with k-fold validation ( $k = 12$ ).

The third model is the one we want to test—this is the one that simulates in basins that are different than the ones that the models were trained on. Out-of-sample testing was done by k-fold validation, which splits the 531 basins randomly into  $k = 12$  groups of approximately equal size, uses all basins from  $k-1$  groups to train the model, and then tests the model on the single group of holdout basins. This procedure is repeated  $k = 12$  times so that out-of-sample predictions are available from every basin. The second model sets an upper benchmark for our PUB LSTMs. In particular, comparison between the second and third models tells us how much information was lost due to prediction in out-of-sample basins versus in-sample basins. Similarly, a comparison between the first and second models lets us evaluate the value of adding catchment attributes to the model inputs, since these are what will, at least potentially, allow the model to be transferable between catchments.

For each model type we trained and tested an ensemble of  $N = 10$  LSTM models to match the 10 SCE restarts used to calibrate the SAC-SMA models. All metrics reported in Section 4 were calculated from the mean of the 10-member ensembles, except for the NWM reanalysis.

All LSTM models were trained on the first 15 years of CAMELS data (1981–1995 water years)—this is the same data period that Newman et al. (2015) used to calibrate SAC-SMA. And all models (LSTMs, SAC-SMA, and NWM) were evaluated on the last 15 years of CAMELS data (1996–2010 water years). LSTMs were trained and evaluated using a k-fold approach ( $k = 12$ ). The training loss function was the average NSE over all training catchments; this is a squared-error loss function that, unlike a more traditional MSE loss function, does not overweight catchments with larger mean streamflow values (i.e., does not overweight large, humid catchments) (Kratzert et al., 2019).

## 4. Results

A comparison between interpolated frequency distributions over the NSE values from 531 CAMELS catchments from all three LSTM models and both benchmark models (SAC-SMA, NWM) is shown in Figure 2.

**Table 2**  
*Summary of Benchmark Statistics for All Models Across 531 Catchments*

	Median	Mean	Minimum	Maximum
Nash Sutcliffe efficiency:	( $-\infty, 1]$ – values close to 1 are desirable.			
SAC-SMA:	0.64	0.51	-12.28	0.88
NWM:	0.58	0.31	-20.28	0.89
Global LSTM (no statics):	0.63	0.45	-31.72	0.90
Global LSTM (with statics):	0.74	0.68	-1.78	0.93
PUB LSTM:	0.69	0.54	-13.02	0.90
Fractional Bias:	( $-\infty, 1]$ – values close to 0 are desirable.			
SAC-SMA:	0.04	0.02	-1.76	0.71
NWM:	0.05	-0.01	-4.80	1.00 <sup>d</sup>
Global LSTM (no statics):	0.01	-0.03	-3.01	0.77
Global LSTM (with statics):	-0.01	-0.04	-2.19	0.49
PUB LSTM:	-0.02	-0.09	-4.86	0.72
Standard Deviation Ratio <sup>e</sup> :	[0, $\infty$ ) – values close to 1 are desirable.			
SAC-SMA:	0.83	0.87	0.10	3.76
NWM:	0.86	0.93	0.00 <sup>f</sup>	4.04
Global LSTM (no statics):	0.74	0.81	0.10	5.83
Global LSTM (with statics):	0.88	0.89	0.17	1.96
PUB LSTM:	0.86	0.91	0.10	3.23
95th Percentile Difference <sup>b</sup> :	( $-\infty, 1]$ – values close to 0 are desirable.			
SAC-SMA:	0.02	-0.05	-3.98	0.83
NWM:	0.07	-0.07	-8.59	1.00 <sup>c</sup>
Global LSTM (no statics):	0.12	0.02	-4.97	0.81
Global LSTM (with statics):	0.03	-0.03	-3.30	0.63
PUB LSTM:	0.03	-0.08	-5.26	0.78

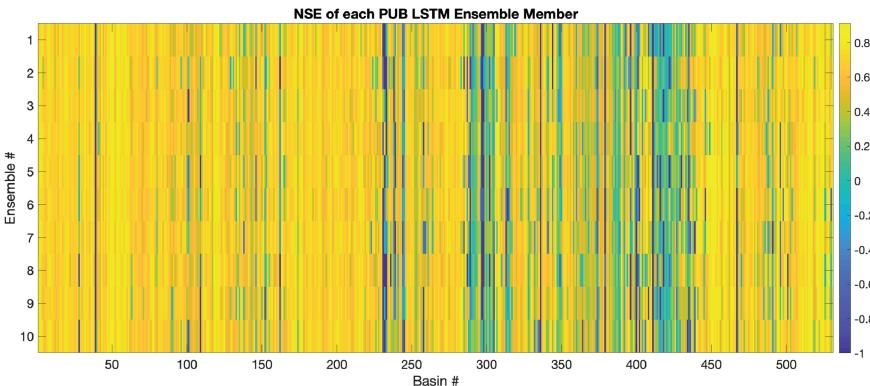
<sup>a</sup>Ratio of the standard deviation of simulated versus observed flows at each catchment. <sup>b</sup>Difference between the values of the observed versus simulated 95th percentile flows divided by the observed 95th percentile flows at each catchment.

<sup>c</sup>Values of zero and one in the NWM max/min statistics are due to rounding. In particular, for one basin (USGS basin ID: 2108000) the NWM simulates a 95th flow percentile of  $\sim 1 \times 10^{-3}$  (mm/day), whereas the 95th percentile of observed flow is  $\sim 4$  (mm/day).

Mean and median values of several performance statistics are given in Table 2. Interpolation was done with kernel density estimation using Gaussian kernels and an optimized bandwidth.

The primary result is that the out-of-sample PUB LSTM ensemble performed at least as well as both of the in-sample benchmarks in more than half of the catchments against all four performance metrics that we tested, except that the basin-calibrated SAC-SMA has a slightly lower average difference between the 95th percentile flows (both SAC-SMA and the PUB LSTM underestimated peak flows to some extent). The PUB LSTM had a higher NSE than SAC-SMA in 307 of 531 (58%) catchments, and higher than the NWM in 347 of 531 (66%) catchments. The PUB LSTM ensemble also had higher mean and maximum NSE scores than the benchmark models; however, SAC-SMA tended to outperform the PUB LSTM in catchments with low NSE values (see the CDF plot in Figure 2).

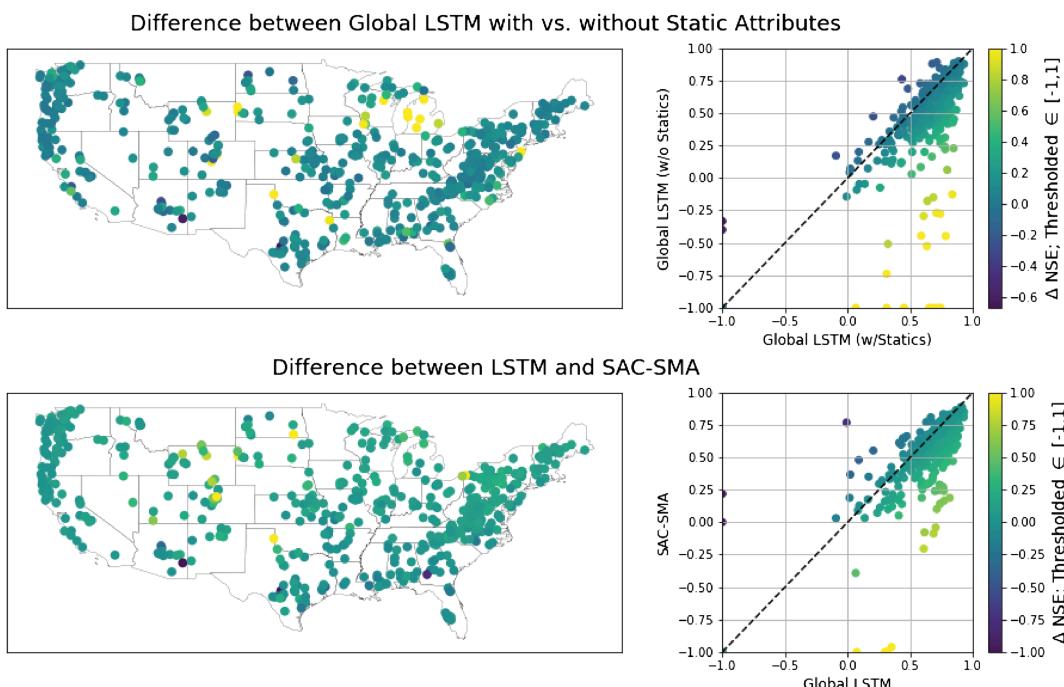
There is some amount of stochasticity associated with training the LSTMs, especially through the random weight initialization of the LSTMs, but also by the weight optimization strategy (we used an ADAM optimizer, Kingma & Ba, 2014). Because of this, the LSTM-type models give better predictions when used as an ensemble. It is not necessarily the case that if one particular LSTM model performs poorly in one catchment that a different LSTM trained on *exactly the same data* will also perform poorly. In our case, we used an ensemble of  $N = 10$  (the same size as the SAC-SMA ensemble developed by Newman et al., 2015 that was used here for benchmarking). Figure 3 shows the NSE values for each ensemble member for the PUB LSTM models. In total, there were 103 basins with at least one PUB LSTM ensemble member with an NSE



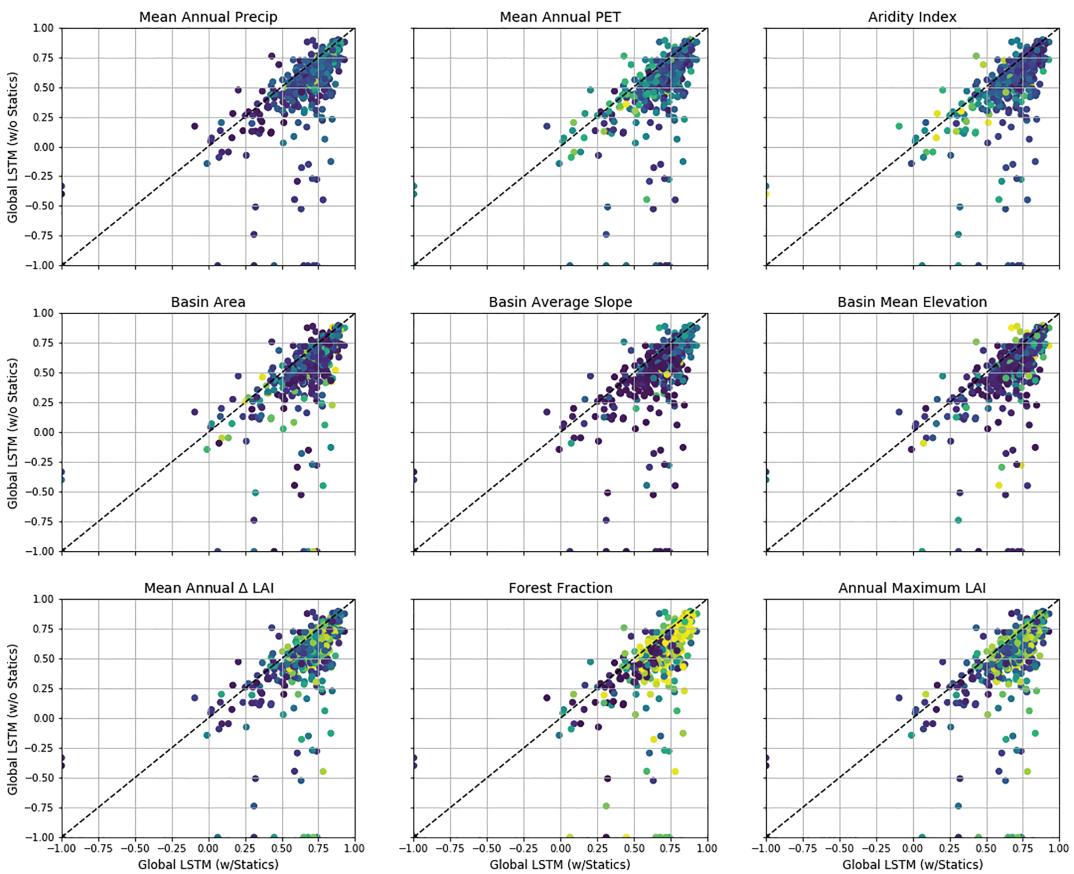
**Figure 3.** NSE scores for all PUB LSTM ensemble members. In some number of basins, certain ensemble members perform well and certain ensemble members perform poorly. This motivates the use of ensembles of LSTMs.

score of below zero. Only 9 of these 103 basins have all  $N = 10$  ensemble members with  $\text{NSE} < 0$ , while 55 of the 103 have at least one ensemble member with  $\text{NSE} > 0.5$ . As an example, one of the basins (USGS basin ID: 01142500, which is basin number 232 in Figure 3) had 9 of 10 ensemble members with  $\text{NSE} < 0$ , but one ensemble member with  $\text{NSE} > 0.7$ . This indicates that a substantial portion of the uncertainty in these LSTM models is due to randomness, rather than to systematic model structural error.

The global LSTM model with static catchment attributes performs better than all other models against the metrics that we tested. Figure 4 compares the performance of the Global LSTM with other benchmark models (SAC-SMA and the Global LSTM without static catchment attributes). The Global LSTM with catchment attributes performs better in most—but not all—catchments. This indicates two things. First, the comparison between the Global LSTM with and without static catchment attributes indicates that although there is useful information in the catchment attributes, in some catchments having these data actually hurts us. We explored this relationship briefly, but did not find any patterns in terms of which catchment attributes might tend to lead to underperformance. Specifically, Figure 5 shows that there is generally no correlation between



**Figure 4.** Comparison between the Global LSTM model with static catchment attributes and other benchmark models used in this study: (top row) LSTM without static catchment attributes and (bottom row) SAC-SMA.



**Figure 5.** Scatterplots of the LSTM NSE scores in each basin with versus without static catchment attributes as model inputs. Colors indicate the relative values of a subselection of the static catchment attributes from Table 1—each subplot has a different colorscale depending on the absolute magnitudes of the specific attributes data. It is the relative values of the attributes that we care about here. There are no apparent direct relationships between the values of different catchment attributes and basins where adding catchment attribute data hurts model performance.

the values of individual catchment attributes and whether the Global LSTM with versus without statics performs better. Our initial conclusion is that the basins where the LSTM without catchment attributes performs better is likely an indication of error or uncertainty in the catchment attributes data. Nonetheless, these data did generally add significant skill to the model (the difference in NSE scores was statistically different at  $p < 1e-9$ ). Future work might use a more sophisticated sensitivity analysis (e.g., sequential model building or a Sobol'-type analysis) to test which specific catchment attributes cause this underperformance when added to the model.

The second thing that we want to highlight from the comparison between the Global LSTM and SAC-SMA (Figure 4) is that there is substantial room to improve SAC-SMA overall. This clearly shows that the LSTM finds rainfall-runoff relationships in individual catchments that SAC-SMA cannot emulate. However, the fact that SAC-SMA performs better in some catchments indicates the potential value of having physical constraints in a hydrological model. The LSTMs in these cases are either overfit or are not able to simulate behaviors of certain similar catchments in the training data set.

## 5. Discussion

The results illustrated in the previous section tell us three things:

1. The process-driven hydrology models that we used here as benchmarks could be improved. The LSTM often finds a better functional representation of rainfall-runoff behavior in most catchments than either SAC-SMA or the NWM.

2. The argument that process-driven models may be preferable in out-of-sample conditions may not hold water. Modern ML methods are quite powerful at extracting information from large, diverse data sets under a variety of hydrological conditions.
3. The comparison between models with and without static catchment attributes as inputs demonstrates that there is sufficient information contained in catchment attribute data to distinguish between different rainfall-runoff relationships in at least most of the U.S. catchments that we tested.

Related to the third conclusion, the challenge going forward is about how to extract the useful information from catchment attributes data for regional modeling. One of the historical reasons why this has been a hard problem is because the usual strategy is to use observable catchment attributes or characteristics to identify or “regionalize” parameters of conceptual or process-based simulation models (e.g., Prieto et al., 2019; Razavi & Coulibaly, 2012). This is hard because of strong interactions in high-dimensional parameter spaces. There are many methods for this—notably a family of regionalization methods that Razavi and Coulibaly (2012) called “model independent”; however, we are unaware of any approach that is as effective as LSTMs at extracting this information for streamflow simulation. This is also in line with the recent results by Kratzert et al. (2019), where similar LSTMs were compared against models calibrated with a parameter regionalization strategy (Samaniego et al., 2010). That paper additionally showed that the response of LSTM-type models were relatively smooth with respect to perturbing catchment attributes, indicating a robust fit (i.e., that the models were not overfit or simply remembering different catchments). The results presented here show that the LSTM is able to extrapolate on catchment attributes to new catchments. Taken together, these results indicate that the catchment attribute data set (Addor et al., 2017) contains a significant amount of useful information about the differences between rainfall-runoff behaviors across (eco)hydrological regimes, and that machine learning is effective at extracting and using these patterns.

Related to the first conclusion, this is yet another example where traditional hydrological models do not take full advantage of the information available from the Earth-observation data record. In this case, neither SAC-SMA nor the NWM are able to directly use the catchment attribute data that we use here, but even if those models could leverage this information, they still could not compete with the LSTM, since the LSTM outperforms even when the conceptual model is calibrated in-basin. This means that not only is there useful information in catchment attributes data, but *also* that there is more information in meteorological forcing data than is used by the traditional models. Several recent experiments have shown the same thing for a number of operational terrestrial hydrology models (e.g., Nearing et al., 2018, 2016). Hrachowitz et al. (2013) and others suggest that better process-based understanding of catchment behaviors should result in better out-of-sample predictions. In reality, it is data-driven models that have consistently given increasingly better predictions. From a more optimistic perspective, ML benchmarking experiments like the one in this paper show that there are probably organizing theories about watersheds yet to be discovered, since machine learning models are able to find informative patterns in multibasin data sets that our current models do not reproduce.

The power of big data and machine learning for problems like this is that such techniques can synthesize information from multiple sites and situations into a single model. As an example, if we were to want to simulate catchment behavior under nonstationary conditions (e.g., evolving climate), then a single LSTM trained to recognize and distinguish different types of hydrological behavior (as shown here) will have a larger range of conditions where it can be expected to remain realistic than a model calibrated to a past conditions in only a single basin.

In our opinion, the most effective strategy moving forward will probably be theory-guided data-science Karpatne et al. (2017). There are now numerous strategies across scientific disciplines that allow for meaningful fusions of domain knowledge with machine learning and other algorithms for learning and predicting directly from data. Adopting approaches like this will be critical moving forward.

## 6. Code and Data Availability

CAMELS data, including SAC-SMA simulations, are available from NCAR at this site (<https://ral.ucar.edu/solutions/products/camels>). National Water Model reanalysis data are available from the NOAA Big Data Repository (<https://registry.opendata.aws/nwm-archive/>). All code used for this project is available at this site ([https://github.com/kratzert/lstm\\_for\\_pub](https://github.com/kratzert/lstm_for_pub)).

## Acknowledgments

The project relies heavily on open source software. All programming was done in Python version 3.7 (van Rossum, 1995) and associated libraries including: Numpy (Van Der Walt et al., 2011), Pandas (McKinney, 2010), PyTorch (Paszke et al., 2017), and Matplotlib (Hunter, 2007). This work was supported by Bosch, ZF, and Google. We thank the NVIDIA Corporation for the GPU donations, LIT with Grant LIT-2017-3-YOU-003 and FWF Grant P 28660-N31.

## References

- Addor, N., Nearing, G., Prieto, C., Newman, A., Le Vine, N., & Clark, M. (2018). A ranking of hydrological signatures based on their predictability in space. *Water Resources Research*, 54, 8792–8812. <https://doi.org/10.1029/2018WR022606>
- Addor, N., Newman, A., Mizukami, N., & Clark, M. P. (2017). Catchment attributes for large-sample studies. <https://doi.org/10.5065/D6G73C3Q>
- Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017). The CAMELS data set: Catchment attributes and meteorology for large-sample studies. *Hydrology and Earth System Sciences (HESS)*, 21(10), 5293–5313.
- AlQuraishi, M. (2019). End-to-end differentiable learning of protein structure. *Cell systems*, 8(4), 292–301.
- Blöschl, G. (2016). Predictions in ungauged basins—Where do we stand? *Proceedings of the International Association of Hydrological Sciences*, 373, 57–60.
- Duan, Q., Gupta, V. K., & Sorooshian, S. (1993). Shuffled complex evolution approach for effective and efficient global minimization. *Journal of optimization theory and applications*, 76(3), 501–521.
- Fekete, B. M., Robarts, R. D., Kumagai, M., Nachtnebel, H.-P., Odada, E., & Zhulidov, A. V. (2015). Time for in situ renaissance. *Science*, 349(6249), 685–686.
- Goswami, M., Oconnor, K., & Bhattacharai, K. (2007). Development of regionalisation procedures using a multi-model approach for flow simulation in an ungauged catchment. *Journal of Hydrology*, 333(2–4), 517–531.
- He, S., Li, Y., Feng, Y., Ho, S., Ravanbakhsh, S., Chen, W., & Póczos, B. (2019). Learning to predict the cosmological structure formation. *Proceedings of the National Academy of Sciences*, 116, 13,825–13,832.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735–1780.
- HRachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivapalan, M., Pomeroy, J., et al. (2013). A decade of Predictions in Ungauged Basins (PUB)—A review. *Hydrological sciences journal*, 58(6), 1198–1255.
- Hsu, K.-I., Gupta, H. V., & Sorooshian, S. (1995). Artificial neural network modeling of the rainfall-runoff process. *Water resources research*, 31(10), 2517–2530.
- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing In Science & Engineering*, 9(3), 90–95.
- Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., et al. (2017). Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Transactions on Knowledge and Data Engineering*, 29(10), 2318–2331.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Kirchner, J. W. (2006). Getting the right answers for the right reasons: Linking measurements, analyses, and models to advance the science of hydrology. *Water Resources Research*, 42, W03S04. <https://doi.org/10.1029/2005WR005362>
- Kleméš, V. (1986). Dilettantism in hydrology: Transition or destiny? *Water Resources Research*, 22(9S), 177S–188S.
- Kratzert, F., Herrnegger, M., Klotz, D., Hochreiter, S., & Klambauer, G. (2018). Do internals of neural networks make sense in the context of hydrology? In *Proceedings of the 2018 AGU fall meeting*. Washington, DC.
- Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall-runoff modelling using long short-term memory (LSTM) networks. *Hydrology and Earth System Sciences*, 22(11), 6005–6022.
- Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019). Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets. *Hydrology and Earth System Sciences*, 23, 5089–5110. <https://doi.org/10.5194/hess-23-5089-2019>
- Liu, Y., Racah, E., Correa, J., Khosrowshahi, A., Lavers, D., Kunkel, K., et al. et al. (2016). Application of deep convolutional neural networks for detecting extreme weather in climate datasets. arXiv preprint arXiv:1605.01156.
- Mayr, A., Klambauer, G., Unterthiner, T., & Hochreiter, S. (2016). Deeptox: Toxicity prediction using deep learning. *Frontiers in Environmental Science*, 3, 80.
- McAfee, A., & Brynjolfsson, E. (2017). *Machine, platform, crowd: Harnessing our digital future*. New York, NY: WW Norton & Company.
- McKinney, W. (2010). Data structures for statistical computing in Python. *Proceedings of the 9th Python in Science Conference*, 1697900(Scipy), 51–56.
- Milly, P. C. D., Betancourt, J., Falkenmark, M., Hirsch, R. M., Kundzewicz, Z. W., Lettenmaier, D. P., & Stouffer, R. J. (2008). Stationarity is dead: Whither water management? *Science*, 319(5863), 573–574.
- Nearing, G. S., & Gupta, H. V. (2015). The quantity and quality of information in hydrologic models. *Water Resources Research*, 51, 524–538. <https://doi.org/10.1002/2014WR015895>
- Nearing, G. S., Mocko, D. M., Peters-Lidard, C. D., Kumar, S. V., & Xia, Y. (2016). Benchmarking NLDAS-2 soil moisture and evapotranspiration to separate uncertainty contributions. *Journal of Hydrometeorology*, 17(3), 745–759.
- Nearing, G. S., Ruddell, B. L., Clark, M. P., Nijssen, B., & Peters-Lidard, C. (2018). Benchmarking and process diagnostics of land models. *Journal of Hydrometeorology*, 19(11), 1835–1852.
- Newman, A., Clark, M., Sampson, K., Wood, A., Hay, L., Bock, A., et al. (2015). Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: Data set characteristics and assessment of regional variability in hydrologic model performance. *Hydrology and Earth System Sciences*, 19(1), 209–223.
- Newman, A., Sampson, K., Clark, M. P., Bock, A., Viger, R. J., & Blodgett, D. (2014). A large-sample watershed-scale hydrometeorological dataset for the contiguous USA. Boulder, CO: UCAR/NCAR. <https://doi.org/10.5065/D6MW2F4D>
- Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., & Nearing, G. (2017). Benchmarking of a physically based hydrologic model. *Journal of Hydrometeorology*, 18(8), 2215–2225.
- Parajka, J., Viglione, A., Rogger, M., Salinas, J., Sivapalan, M., & Blöschl, G. (2013). Comparative assessment of predictions in ungauged basins—Part 1: Runoff-hydrograph studies. *Hydrology and Earth System Sciences*, 17(5), 1783–1795.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., et al. (2017). Automatic differentiation in PyTorch.
- Prieto, C., Le Vine, N., Kavetski, D., Garcia, E., & Medina, R. (2019). Flow prediction in ungauged catchments using probabilistic random forests regionalization and new statistical adequacy tests. *Water Resources Research*, 55, 4364–4392. <https://doi.org/10.1029/2018WR023254>
- Razavi, T., & Coulibaly, P. (2012). Streamflow prediction in ungauged basins: Review of regionalization methods. *Journal of Hydrologic Engineering*, 18(8), 958–975.
- Salas, F. R., Somos-Valenzuela, M. A., Dugger, A., Maidment, D. R., Gochis, D. J., David, C. H., et al. (2018). Towards real-time continental scale streamflow simulation in continuous and discrete space. *JAWRA Journal of the American Water Resources Association*, 54(1), 7–27.
- Samaniego, L., Kumar, R., & Attinger, S. (2010). Multiscale parameter regionalization of a grid-based hydrologic model at the mesoscale. *Water Resources Research*, 46, W05523. <https://doi.org/10.1029/2008WR007327>

- Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., et al. (2017). Toward seamless hydrologic predictions across spatial scales. *Hydrology and Earth System Sciences*, 21(9), 4323–4346. <https://doi.org/10.5194/hess-21-4323-2017>
- Sellars, S. (2018). “Grand challenges” in big data and the earth sciences. *Bulletin of the American Meteorological Society*, 99(6), ES95–ES98.
- Sivapalan, M. (2003). Prediction in ungauged basins: A grand challenge for theoretical hydrology. *Hydrological Processes*, 17(15), 3163–3170.
- Sivapalan, M., Takeuchi, K., Franks, S., Gupta, V., Karambiri, H., Lakshmi, V., et al. (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.
- Van Der Walt, S., Colbert, S. C., & Varoquaux, G. (2011). The NumPy array: A structure for efficient numerical computation. *Computing in Science and Engineering*, 13(2), 22–30.
- van Rossum, G. (1995). Python tutorial (*Technical Report CS-R9526*). Amsterdam: Centrum voor Wiskunde en Informatica (CWI).
- Vaze, J., Chiew, F., Hughes, D., & Andréassian, V. (2015). Preface: Hs02—hydrologic non-stationarity and extrapolating models to predict the future. *Proceedings of the International Association of Hydrological Sciences*, 371, 1–2.
- Vrugt, J. A., Gupta, H. V., Dekker, S. C., Sorooshian, S., Wagener, T., & Bouten, W. (2006). Application of stochastic parameter optimization to the Sacramento Soil Moisture Accounting Model. *Journal of Hydrology*, 325(1-4), 288–307.
- Xia, Y., Mitchell, K., Ek, M., Cosgrove, B., Sheffield, J., Luo, L., et al. (2012). Continental-scale water and energy flux analysis and validation for North American land data assimilation system project phase 2 (NLDAS-2): 2. Validation of model-simulated streamflow. *Journal of Geophysical Research*, 117, D03110. <https://doi.org/10.1029/2011JD016051>