

[Related articles](#)
[HESS](#) | [Articles](#) | [Volume 22, issue 11](#)

Hydrol. Earth Syst. Sci., 22, 6005–6022, 2018

<https://doi.org/10.5194/hess-22-6005-2018>

© Author(s) 2018. This work is distributed under the Creative Commons Attribution 4.0 License.



Research article | 22 Nov 2018

# Rainfall–runoff modelling using Long Short-Term Memory (LSTM)

**Frederik Kratzert**<sup>ID1,\*</sup>, **Daniel Klotz**<sup>ID1</sup>, **Claire Brenner**<sup>ID1</sup>, **Karsten Schulz**<sup>ID1</sup>, and **Mathew Herrnegger**<sup>ID1</sup>
<sup>1</sup>Institute of Water Management, Hydrology and Hydraulic Engineering, University of Natural Resources and Life Sciences, Vienna, 1190, Austria

\* Invited contribution by Frederik Kratzert, recipient of the EGU Hydrological Sciences Outstanding Student Poster and PICO Award 2016.

**Correspondence:** Frederik Kratzert (f.kratzert@gmail.com)**Received: 04 May 2018 – Discussion started: 14 May 2018 – Revised: 25 Sep 2018 – Accepted: 13 Nov 2018 – Published: 22 Nov 2018**

## Abstract

Rainfall–runoff modelling is one of the key challenges in the field of hydrology. Various approaches exist, ranging from physically based conceptual to fully data-driven models. In this paper, we propose a novel data-driven approach, using the Long Short-Term Memory (LSTM) network, a special type of recurrent neural network. The advantage of the LSTM is its ability to learn long-term dependencies between provided input and output of the network, which are essential for modelling storage effects in e.g. catchments with snow influence. We apply the LSTM to 241 catchments of the freely available CAMELS data set to test our approach and also compare the results to the well-known Sacramento Soil Moisture Accounting Model (SAC-SMA) coupled with the Snow-17 snow routine. We also show the potential of the LSTM as a regional hydrological model in which one model predicts the discharge for a variety of catchments. In our last experiment, we show the potential of transfer process understanding, learned at regional scale, to individual catchments and thereby increasing model performance when compared to a LSTM trained only on the data of single catchments. Using this approach, we were able to achieve better model performance than the SAC-SMA + Snow-17, which underlines the potential of the LSTM for hydrological modelling applications.

### Download & links

[Article \(PDF, 10256 KB\)](#)

**How to cite.** Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall–runoff modelling using Long Short-Term Memory (LSTM) networks, Hydrol. Earth Syst. Sci., 22, 6005–6022, <https://doi.org/10.5194/hess-22-6005-2018>, 2018.

## 1 Introduction

Rainfall–runoff modelling has a long history in hydrological sciences and the first attempts to predict the discharge as a function of precipitation events using regression-type approaches date back 170 years (Beven, 2001; Mulvaney, 1850). Since then, modelling approaches have been further developed by progressively incorporating physically based process understanding and concepts into the (mathematical) model formulations. These include explicitly addressing the spatial variability of processes, boundary conditions and physical processes in catchments (Freeze and Harlan, 1969; Kirchner, 2006; Schulla, 2007). These developments are largely driven by the advancement of computer technology and the availability of (remote sensing) data at high spatial and temporal resolution (Hengl et al., 2017; Kratzert et al., 2010; Mu et al., 2011; Myneni et al., 2002; Rennó et al., 2008).

However, the development towards coupled, physically based and spatially explicit representations of hydrological processes at the catchment scale has come at the price of high computational costs and a high demand for necessary (meteorological) input data (Wood et al., 2016). Therefore, physically based models are still rarely used in operational rainfall–runoff forecasting. In addition, the current data scarcity and parameterization of these kind of models, e.g. the 3-D information on the physical characteristics of the sub-surface, are mostly limited to small, experimental watersheds, limiting the model's applicability for larger river basins in an operational context. The high computational costs further limit their application, especially if uncertainty estimations and multiple model runs within an ensemble forecasting

hydrological problems. [Shi et al. \(2015\)](#) investigated a deep learning approach for precipitation nowcasting. [Tao et al. \(2016\)](#) used a neural network for bias correction of satellite precipitation products. [Fang et al. \(2017\)](#) investigated the use of deep learning models for soil moisture in the context of NASA's Soil Moisture Active Passive (SMAP) satellite mission. [Assem et al. \(2017\)](#) compared the performance of a deep learning approach for water flow level and flow predictions for the Shannon River in Ireland with multiple baseline models and reported that the deep learning approach outperforms all baseline models consistently. More recently, [D. Zhang et al. \(2018\)](#) compared the performance of different neural network architectures for simulating and predicting the water levels of a combined sewer structure in Drammen (Norway), based on online data from rain gauges and water-level sensors. They confirmed that LSTM (as well as another neural network architecture with cell memory) are better suited for multi-step-ahead predictions than traditional architectures without explicit cell memory. [J. Zhang et al. \(2018\)](#) used an LSTM for predicting water tables in agricultural areas. Among other things, they compared the resulting simulation from the LSTM-based approach with that of a traditional neural network and found that the LSTM outperforms the latter. In general, the potential use and benefits of DL approaches in the field of hydrology and water sciences have recently come into the focus of discussion ([Marçais and de Dreuzy, 2017](#); [Shen, 2018](#); [Shen et al., 2018](#)). In this context we would like to mention [Shen \(2018\)](#) more explicitly, since he provides an ambitious argument for the potential of DL in earth sciences/hydrology, so he also provides an overview of various applications of DL in earth sciences. Of special interest for the present case is his point that DL might also provide an avenue for discovering emergent behaviours of hydrological phenomena.

Regardless of the hydrological modelling approach applied, any model will be typically calibrated for specific catchments for which a time series of meteorological and hydrological data are available. The calibration procedure is required because models are only simplifications of real catchment hydrology and model parameters have to effectively represent non-resolved processes and any subgrid-scale heterogeneity in catchment characteristics (e.g. soil hydraulic properties) ([Beven, 1995](#); [Merz et al., 2006](#)). The transfer of model parameters (regionalization) from catchments where meteorological and runoff data are available to ungauged or data-poor catchments is one of the ongoing challenges in hydrology ([Buytaert and Beven, 2009](#); [He et al., 2011](#); [Samaniego et al., 2010](#)).

The aim of this study is to explore the potential of the LSTM architecture (in the adapted version proposed by [Gers et al., 2000](#)) for the rainfall–runoff behaviour of a large number of differently complex catchments at the daily timescale. Additionally, we want to explore the potential of LSTMs for regionalizing the rainfall–runoff response by training a single model for a multitude of catchments. In order to draw a more general conclusion about the suitability of our modelling approach, we test this approach on a large number of catchments from the CAMELS data set ([Addor et al., 2017b](#); [Newman et al., 2014](#)). This data set is freely available and includes meteorological forcing and observed discharge for 671 catchments across the contiguous United States. For each basin, the CAMELS data set also includes the simulated discharge from the Sacramento Soil Moisture Accounting Model ([Burnash et al., 1973](#)) coupled with the Snow-17 snow storage model ([Anderson, 1973](#)). In our study, we use these simulations as a benchmark, to compare our model results with an established modelling approach.

The paper is structured in the following way: in Sect. 2, we will briefly describe the LSTM network architecture and the data set used, followed by an introduction into three different experiments: in the first experiment, we test the general ability of the LSTM to model runoff processes for a large number of individual catchments. The second experiment investigates the capability of LSTMs for regional modelling, and the last tests whether the regional models can help to enhance the simulation performance for individual catchments. Section 3 presents and discusses the results of our experiments, before we end our paper with a conclusion and outlook for future work.

## 2 Methods and database

### 2.1 Long Short-Term Memory network

In this section, we introduce the LSTM architecture in more detail, using the notation of [Graves et al. \(2013\)](#). Beside a technical description of the network internals, we added a “hydrological interpretation of the LSTM” in Sect. 3.5 in order to bridge differences between the hydrological and deep learning research communities.

The LSTM architecture is a special kind of recurrent neural network (RNN), designed to overcome the weakness of the traditional RNN to handle long-term dependencies. [Bengio et al. \(1994\)](#) have shown that the traditional RNN can hardly remember sequences with a length greater than 10. For daily streamflow modelling, this would imply that we could only use the last 10 days of meteorological data as input to predict the streamflow of the next day. This period is too short considering the memory of catchments including groundwater, snow or even surface water storages, with lag times between precipitation and discharge up to several years.

To explain how the RNN and the LSTM work, we unfold the recurrence of the network into a directed acyclic graph (see Fig. 1). The output (in our case discharge) for a specific time step is predicted from the input  $x=[x_1, \dots, x_n]$  consisting of the last  $n$  consecutive time steps. The input consists of independent variables (in our case daily precipitation, min/max temperature, solar radiation and vapour pressure) and is processed sequentially. In each time step  $t$  ( $1 \leq t \leq n$ ), the current input  $x_t$  is processed in the recurrent cells of each layer in the network.

**Figure 2 (a)** The internal operation of a traditional RNN cell:  $\mathbf{h}_t$  stands for hidden state and  $\mathbf{x}_t$  for the input at time step  $t$ . **(b)** internals of a LSTM cell, where  $\mathbf{f}$  stands for the forget gate (Eq. 2),  $\mathbf{i}$  for the input gate (Eqs. 3–4), and  $\mathbf{o}$  for the output gate (Eq. 5).  $\mathbf{c}_t$  denotes the cell state at time step  $t$  and  $\mathbf{h}_t$  the hidden state.

► Download

In a traditional RNN cell, only one internal state  $\mathbf{h}_t$  exists (see Fig. 2a), which is recomputed in every time step by the following equation:

$$\mathbf{h}_t = g(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}),$$

where  $g(\cdot)$  is the activation function (typically the hyperbolic tangent),  $\mathbf{W}$  and  $\mathbf{U}$  are the adjustable weight matrices of the hidden state and the input  $\mathbf{x}$ , and  $\mathbf{b}$  is an adjustable bias vector. In the first time step, the hidden state is initialized as a vector of zeros and its length is a user-defined hyperparameter of the network.

In comparison, the LSTM has (i) an additional cell state or cell memory  $\mathbf{c}_t$  in which information can be stored, and (ii) gates (three letters in Fig. 2b) that control the information flow within the LSTM cell (Hochreiter and Schmidhuber, 1997). The first gate is the forget gate, introduced by Gers et al. (2000). It controls which elements of the cell state vector  $\mathbf{c}_{t-1}$  will be forgotten (to which degree):

$$\mathbf{f}_t = \sigma(\mathbf{W}_f\mathbf{x}_t + \mathbf{U}_f\mathbf{h}_{t-1} + \mathbf{b}_f),$$

where  $\mathbf{f}_t$  is a resulting vector with values in the range (0, 1),  $\sigma(\cdot)$  represents the logistic sigmoid function and  $\mathbf{W}_f$ ,  $\mathbf{U}_f$  and  $\mathbf{b}_f$  define learnable parameters for the forget gate, i.e. two adjustable weight matrices and a bias vector. As for the traditional RNN, the hidden state is initialized in the first time step by a vector of zeros with a user-defined length.

In the next step, a potential update vector for the cell state is computed from the current input ( $\mathbf{x}_t$ ) and the last hidden state ( $\mathbf{h}_{t-1}$ ) by the following equation:

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_{\tilde{c}}\mathbf{x}_t + \mathbf{U}_{\tilde{c}}\mathbf{h}_{t-1} + \mathbf{b}_{\tilde{c}}),$$

where  $\tilde{\mathbf{c}}_t$  is a vector with values in the range (-1, 1),  $\tanh(\cdot)$  is the hyperbolic tangent and  $\mathbf{W}_{\tilde{c}}$ ,  $\mathbf{U}_{\tilde{c}}$  and  $\mathbf{b}_{\tilde{c}}$  are another set of learnable parameters.

Additionally, the second gate is computed, the input gate, defining which (and to what degree) information of  $\tilde{\mathbf{c}}_t$  is used to update the cell state in the current time step:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i\mathbf{x}_t + \mathbf{U}_i\mathbf{h}_{t-1} + \mathbf{b}_i),$$

where  $\mathbf{i}_t$  is a vector with values in the range (0, 1), and  $\mathbf{W}_i$ ,  $\mathbf{U}_i$  and  $\mathbf{b}_i$  are a set of learnable parameters, defined for the input gate.

With the results of Eqs. (2)–(4) the cell state  $\mathbf{c}_t$  is updated by the following equation:

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t,$$

where  $\odot$  denotes element-wise multiplication. Because the vectors  $\mathbf{f}_t$  and  $\mathbf{i}_t$  have both entries in the range (0, 1), Eq. (5) can be interpreted the way that it defines, which information stored in  $\mathbf{c}_{t-1}$  will be forgotten (values of  $\mathbf{f}_t$  of approx. 0) and which will be kept (values of  $\mathbf{f}_t$  of approx. 1). Similarly,  $\mathbf{i}_t$  decides which new information stored in  $\tilde{\mathbf{c}}_t$  will be added to the cell state (values of  $\mathbf{i}_t$  of approx. 1) and which will be ignored (values of  $\mathbf{i}_t$  of approx. 0). Like the hidden state vector, the cell state is initialized by a vector of zeros in the first time step. Its length corresponds to the length of the hidden state vector.

The third and last gate is the output gate, which controls the information of the cell state  $\mathbf{c}_t$  that flows into the new hidden state. The output gate is calculated by the following equation:

$$\mathbf{o}_t = \sigma(\mathbf{W}_o\mathbf{x}_t + \mathbf{U}_o\mathbf{h}_{t-1} + \mathbf{b}_o),$$

where  $\mathbf{o}_t$  is a vector with values in the range (0, 1), and  $\mathbf{W}_o$ ,  $\mathbf{U}_o$  and  $\mathbf{b}_o$  are a set of learnable parameters, defined for the output gate. With this vector, the new hidden state  $\mathbf{h}_t$  is calculated by combining the results of Eqs. (5) and (6):

$$\mathbf{h}_t = \tanh(\mathbf{c}_t) \odot \mathbf{o}_t.$$

It is in particular the cell state ( $\mathbf{c}_t$ ) that allows for an effective learning of long-term dependencies. Due to its very simple linear operation with the remaining LSTM cell, it can store information unchanged over a long period of time steps. During training, this characteristic

```

1Input:  $x = [x_1, \dots, x_n], x_t \in \mathbb{R}^m$ 

2Given parameters:  $W_f, U_f, b_f, W_{\tilde{c}}, U_{\tilde{c}}, b_{\tilde{c}}, W_i, U_i, b_i,$   

 $W_o, U_o, b_o$ 

3Initialize  $h_0, c_0 = \vec{0}$  of length  $p$ 

4for  $t=1, \dots, n$  do

    5Calculate  $f_t$  (Eq. 2),  $\tilde{c}_t$  (Eq. 3),  $i_t$  (Eq. 4)

    6Update cell state  $c_t$  (Eq. 5)

    7Calculate  $o_t$  (Eq. 6),  $h_t$  (Eq. 7)

8end for

9Output:  $h = [h_1, \dots, h_n], h_t \in \mathbb{R}^p$ 

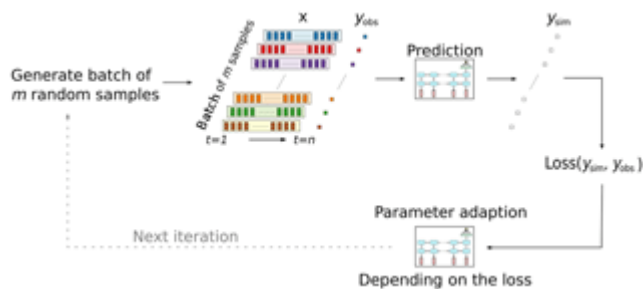
```

## 2.2 The calibration procedure

In traditional hydrological models, the calibration involves a defined number of iteration steps of simulating the entire calibration given set of model parameters and evaluating the model performance with some objective criteria. The model parameters are, in the applied optimization technique (global and/or local), perturbed in such a way that the maximum (or minimum) of an objective function is found. Regarding the training of a LSTM, the adaptable (or *learnable*) parameters of the network, the weights and biases, are also updated depending on a given loss function of an iteration step. In this study we used the mean-squared error (MSE) as an objective criterion.

In contrast to most hydrological models, the neural network exhibits the property of differentiability of the network equations. The gradient of the loss function with respect to any network parameter can always be calculated explicitly. This property is used in the back-propagation step in which the network parameters are adapted to minimize the overall loss. For a detailed description see Goodfellow et al. (2016).

A schematic illustration of one iteration step in the LSTM training/calibration is provided in Fig. 3. One iteration step during the training of LSTMs usually works with a subset (called *batch* or *mini-batch*) of the available training data. The number of samples per batch is a hyperparameter, which in our case was defined to be 512. Each of these samples consists of one discharge value of a given day and the meteorological input of the  $n$  preceding days. In every iteration step, the loss function is calculated as the average of the MSE of the predicted and observed runoff of these 512 samples. Since the discharge of a specific time step is only a function of the meteorological input of the last  $n$  days, the samples within a batch can consist of random time steps (depicted in Fig. 3 by the different colours), which must not necessarily be ordered chronologically. For faster convergence, it is even advantageous to have random samples in one batch (Lecun et al., 2012). This procedure is different from traditional hydrological model calibration, where usually all the information of the training data is processed in each iteration step, since all simulated and observed runoff pairs are used in the model evaluation.



**Figure 3** Illustration of one iteration step in the training process of the LSTM. A random batch of input data  $x$  consisting of  $m$  independent training samples (depicted by the colours) is used in each step. Each training sample consists of  $n$  days of look-back and one target value ( $y_{\text{obs}}$ ) to predict. The loss is computed from the observed discharge and the network's predictions  $y_{\text{sim}}$  and used to update the network parameters.

For efficient learning, all input features (the meteorological variables) as well as the output (the discharge) data are normalized by subtracting the mean and dividing by the standard deviation (LeCun et al., 2012; Minns and Hall, 1996). The mean and standard deviation used for normalization are calculated from the calibration period only. To receive the final discharge prediction, the output of the network is retransformed using the normalization parameters from the calibration period (Fig. 4 shows the retransformed model outputs).

## 2.3 Open-source software

Our research heavily relies on open source software. The programming language of choice is Python 3.6 (van Rossum, 1995). The libraries used for preprocessing our data and for data management in general are Numpy (Van Der Walt et al., 2011), Pandas (McKinney, 2010), Scikit-Learn (Pedregosa et al., 2011). The Deep-Learning frameworks we use are TensorFlow (Abadi et al., 2016) and Keras (Chollet, 2015). All figures are made using Matplotlib (Hunter, 2007).

**Table 1** Overview of the HUCs considered in this study and some region statistics averaged over all basins in that region. For each region, the variable mean and standard deviation is reported.

HUC	Region name	No. of basins	Mean precipitation (mm day <sup>-1</sup> )	Mean aridity <sup>1</sup> (-)	Mean altitude (m)	Mean snow frac. <sup>2</sup> (-)	Mean seasonality <sup>3</sup> (-)
01	New England	27	3.61 ± 0.26	0.60 ± 0.03	316 ± 182	0.24 ± 0.06	0.10 ± 0.08
03	South Atlantic-Gulf	92	3.79 ± 0.49	0.87 ± 0.14	189 ± 179	0.02 ± 0.02	0.12 ± 0.26
11	Arkansas-White-Red	31	2.86 ± 0.89	1.18 ± 0.50	613 ± 713	0.08 ± 0.13	0.25 ± 0.29
17	Pacific Northwest	91	5.22 ± 2.03	0.59 ± 0.40	1077 ± 589	0.33 ± 0.23	-0.72 ± 0.17

<sup>1</sup> PET/P; see Addor et al. (2017a). <sup>2</sup> Fraction of precipitation falling on days with temperatures below 0 °C. <sup>3</sup> Positive values indicate precipitation peaks in summer, negative values that precipitation peaks in the winter month, and values close to 0 that the precipitation is uniform throughout the year (see Addor et al., 2017a).

- Download Print Version
- |
- Download XLSX

## 2.4 The CAMELS data set

The underlying data for our study is the CAMELS data set (Addor et al., 2017b; Newman et al., 2014). The acronym stands for “Catchment Attributes for Large-Sample Studies” and it is a freely available data set of 671 catchments with minimal human disturbances across the contiguous United States (CONUS). The data set contains catchment aggregated (lumped) meteorological forcing data and observed discharge at the daily timescale starting (for most catchments) from 1980. The meteorological data are calculated from three different gridded data sources (Daymet, Thornton et al., 2012; Maurer, Maurer et al., 2002; and NLDAS, Xia et al., 2012) and consists of precipitation, shortwave downward radiation, maximum and minimum temperature, snow-water equivalent and humidity. We use Daymet data, since it has the highest spatial resolution (1 km grid compared to 12 km grid for Maurer and NLDAS) as a basis for the catchment averages and all available meteorological input variables with exception of the snow-water equivalent and the daily streamflow.

The 671 catchments in the data set are grouped into 18 hydrological units (HUCs) following the U.S. Geological Survey's HUC map (Addor et al., 1987). These groups correspond to geographic areas that represent the drainage area of either a major river or the combined area of a series of rivers.

In our study, we used 4 out of the 18 hydrological units with their 241 catchments (see Fig. 5 and Table 1) in order to cover a wide range of different hydrological conditions on one hand and to limit the computational costs on the other hand. The New England region in the east contains 27 more or less homogeneous basins (e.g. in terms of snow influence or aridity). The Arkansas-White-Red region in the south of CONUS has a comparable number of basins, namely 32, but is completely different otherwise. Within this region, attributes e.g. mean annual precipitation have a high variance and strong gradient from east to west (see Fig. 5). Also comparable in size but with different hydro-climatic conditions are the South Atlantic-Gulf region (92 basins) and the Pacific Northwest region (91 basins). The latter stretches from the Pacific coast till the Rocky Mountains and also exhibits a high variance of attributes across the basins, comparable to the Arkansas-White-Red region. For example, there are very humid catchments with more than 3000 mm yr<sup>-1</sup> precipitation close to the Pacific coast (aridity index 2.17, mean annual precipitation 500 mm yr<sup>-1</sup>) basins in the south-east of this region. The relatively flat South Atlantic-Gulf region contains more homogeneous basins, but in contrast to the New England region is not influenced by snow.



## 2.5 Experimental design

Throughout all of our experiments, we used a two-layer LSTM network, with each layer having a cell/hidden state length of 20. Table 2 shows the resulting shapes of all model parameters from Eqs. (2) to (8). Between the layers, we added dropout, a technique to prevent overfitting (Srivastava et al., 2014). Dropout sets a certain percentage (10 % in our case) of random neurons to zero during training in order to force the network into a more robust feature learning. Another hyperparameter is the length of the input sequence, which corresponds to the number of days of meteorological input data provided to the network for the prediction of the next discharge. We decided to keep this value constant at 365 days for this study in order to capture at least the dynamics of a full annual cycle.

The specific design of the network architecture, i.e. the number of layers, cell/hidden state length, dropout rate and input sequence length were found through a number of experiments in several seasonal-influenced catchments in Austria. In these experiments, different architectures (e.g. one or two LSTM layers or 5, 10, 15, or 20 cell/hidden units) were varied manually. The architecture used in this study proved to work well for these catchments (in comparison to a calibrated hydrological model we had available from previous studies, Herrnegger et al., 2018) and was therefore chosen to be applied here without further tuning. A systematic sensitivity analysis of different hyper-parameters was however not done and is something to do in the future.

**Table 2** Shapes of learnable parameters of all layers.

Layer	Parameter	Shape
1st LSTM layer	$W_f, W_{\tilde{c}}, W_i, W_o$	[20, 5]
	$U_f, U_{\tilde{c}}, U_i, U_o$	[20, 20]
	$b_f, b_{\tilde{c}}, b_i, b_o$	[20]
2nd LSTM layer	$W_f, W_{\tilde{c}}, W_i, W_o$	[20, 20]
	$U_f, U_{\tilde{c}}, U_i, U_o$	[20, 20]
	$b_f, b_{\tilde{c}}, b_i, b_o$	[20]
Dense layer	$W_d$	[20, 1]
	$b_d$	[1]

► [Download Print Version](#)

|  
► [Download XLSX](#)

We want to mention here that our calibration scheme (see description in the three experiments below) is not the standard way for training and selecting data-driven models, especially neural networks. As of today, a widespread calibration strategy for DL models is to split the data into three parts, referred to as training, validation and test data (see Goodfellow et al., 2016). The first two splits are used for the parametrization of the networks and the remainder of the data to diagnose the actual performance. We decided to not implement this splitting strategy, because we are limited to the periods Newman et al. (2015) used so that our models are comparable with their results. Theoretically, it would be possible to split the 15-year calibration period of Newman et al. (2015) further into a training and validation period. However, this would lead to (a) a much shorter period of data that is used for the actual weight updates or (b) a high risk of overfitting on a short validation period, depending on how this 15-year period is divided. In addition to that, LSTMs with a low number of hidden units are quite sensitive to the initialization of their weights. It is thus common practice to repeat the calibration task several times with different random seeds to select the best performing realization of the model (Bengio, 2012). For the present purpose we decided not to implement these strategies, since it would make it more difficult or even impossible to compare the LSTM approach to the SAC-SMA + Smoother reference model. The goal of this study is therefore not to find the best per-catchment model, but rather to investigate the general applicability of LSTMs for the task of rainfall–runoff modelling. However, we think that the sample size of 241 catchments is large enough to infer the (average) properties of the LSTM-based approach.

### 2.5.1 Experiment 1: one model for each catchment

With the first experiment, we test the general ability of our LSTM network to model rainfall–runoff processes. Here, we train one model separately for each of the 241 catchments. To avoid the effect of overfitting of the network on the training data, we identified the best model after 100 epochs (for a definition of an epoch, see Sect. 2.2) in a preliminary step, which yielded, on average, the highest Nash–Sutcliffe efficiency (NSE) across all basins for an independent validation period. For this preliminary experiment, we used the first 14 years of the 15-year calibration period as training data and the last, fifteenth, year as the independent validation period. With the 14 years of data, we trained the model for in total 200 epochs for each catchment and evaluated each model after each epoch with the validation data. Across all catchments, the highest mean NSE was achieved after 50 epochs in this preliminary experiment. Thus, for the final training of the LSTM with



same batch size as in Experiment 1 is used (see Sect. 2.2 for an explanation of the connection of number of iterations, number of samples and number of epochs). Thus, for the final training, we train one LSTM for each of the four used HUCs for 20 epochs with 15-year long calibration period.

### 2.5.3 Experiment 3: fine-tuning the regional model for each catchment

In the third experiment, we want to test whether the more general knowledge of the regional model (Experiment 2) can help to performance of the LSTM in a single catchment. In the field of DL this is a common approach called fine-tuning (Razavian et al., Yosinski et al., 2014), where a model is first trained on a huge data set to learn general patterns and relationships between (meteorological) input data and (streamflow) output data (this is referred to as *pre-training*). Then, the pre-trained network is further trained for a number of epochs with the data of a specific catchment alone to adapt the more generally learned processes to a specific catchment. In other speaking, the LSTM first learns the general behaviour of the runoff generating processes from a large data set, and is in a second step adapted in order to account for the specific behaviour of a given catchment (e.g. the scaling of the runoff response in a specific catchment).

In this study, the regional models of Experiment 2 serve as pre-trained models. Therefore, depending on the affiliation of a catchment to a certain HUC, the specific regional model for this HUC is taken as a starting point for the fine-tuning. With the initial LSTM weights from the regional model, the training is continued only with the training data of a specific catchment for a few epochs (ranging from 0 to 20, with a median 10). Thus, similar to Experiment 1, we finally have 241 different models, one for each of the 241 catchments. Different from the previous experiments, we do not use a global number of epochs for fine-tuning. Instead, we used the 14-year/1-year split to determine the optimal number of epochs for each catchment individually. The reason is that the regional model fits individual catchments within the HUC differently well. Therefore, the number of epochs the LSTM needs to adapt to a certain catchment before it starts to overfit is different for each catchment.

## 2.6 Evaluation metrics

The metrics for model evaluation are the Nash–Sutcliffe efficiency (Nash and Sutcliffe, 1970) and the three decompositions following Legates et al. (2009). These are the correlation coefficient of the observed and simulated discharge ( $r$ ), the variance bias ( $\alpha$ ) and the total bias ( $\beta$ ). While all of these measures evaluate the performance over the entire time series, we also use three different signatures of the duration curve (FDC) that evaluate the performance of specific ranges of discharge. Following Yilmaz et al. (2008), we calculate the 2 % flows, the peak flows (FHV), the bias of the slope of the middle section of the FDC (FMS) and the bias of the bottom 30 % (FLV).

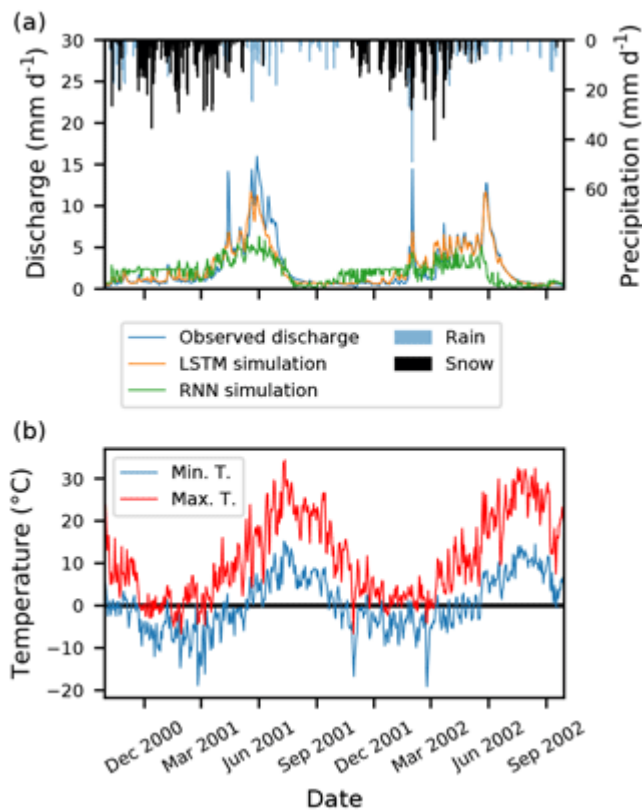
Because our modelling approach needs 365 days of meteorological data as input for predicting one time step of discharge, we calculate the first year of the calibration period. To be able to compare our models to the SAC-SMA + Snow-17 benchmark model, we record the same metrics for the benchmark model for the same simulation periods.

## 3 Results and discussion

We start presenting our results by showing an illustrative comparison of the modelling capabilities of traditional RNNs and the LSTM. We highlight the problems of RNNs to learn long-term dependencies and its deficits for the task of rainfall–runoff modelling. This is followed by the analysis of the results of Experiment 1, for which we trained one network separately for each basin and compare the results to the SAC-SMA + Snow-17 benchmark model. Then we investigate the potential of LSTMs to learn hydrological behaviour at the regional scale. In this context, we compare the performance of the regional models from Experiment 2 against the models of Experiment 1 and discuss their strengths and weaknesses. Lastly, we examine whether our fine-tuning approach enhances the predictive power of our models in individual catchments. In all cases, the analysis is based on the data of the 241 catchments of the calibration (the first 15 years) and validation (all remaining years available) periods.

### 3.1 The effect of (not) learning long-term dependencies

As stated in Sect. 2.1, the traditional RNN can only learn dependencies of 10 or less time steps. The reason for this is the so-called “vanishing or exploding gradients” phenomenon (see Bengio et al., 1994, and Hochreiter and Schmidhuber, 1997), which manifests itself in a signal during the backward pass of the network training that either diminishes towards zero or grows against infinity, preventing the learning of long-term dependencies. However, from the perspective of hydrological modelling, a catchment contains various processes with dependencies well above 10 days (which corresponds to 10 time steps in the case of daily streamflow modelling), e.g. snow accumulation during winter and snowmelt during spring and summer. Traditional hydrological models need to reproduce these processes correctly to be able to make accurate streamflow predictions. This is in principle not the case for data-driven approaches.



**Figure 6 (a)** Two years of observed as well as the simulated discharge of the LSTM and RNN from the validation period of basin 13340600. The precipitation is plotted from top to bottom and days with minimum temperature below zero are marked with black bars. **(b)** The corresponding daily maximum and minimum temperature.

► [Download](#)

In contrast, the LSTM seems to have (i) no or fewer problems with predicting the correct amount of discharge during the snowmelt and (ii) the predicted hydrograph is much smoother and fits the general trends of the hydrograph much better. Note that both models were trained with the exact same data and have the same data available for predicting a single day of discharge.

Here we have only shown a single example for a snow-influenced basin. We also compared the modelling behaviour in one of the catchments of the Arkansas-White-Red region, and found that the trends and conclusion were similar. Although only based on an illustrative example that shows the problems of RNNs with long-term dependencies, we can conclude that traditional RNNs should not be used if (e.g. daily) discharge is predicted only from meteorological observations.

## 3.2 Using LSTMs as hydrological models

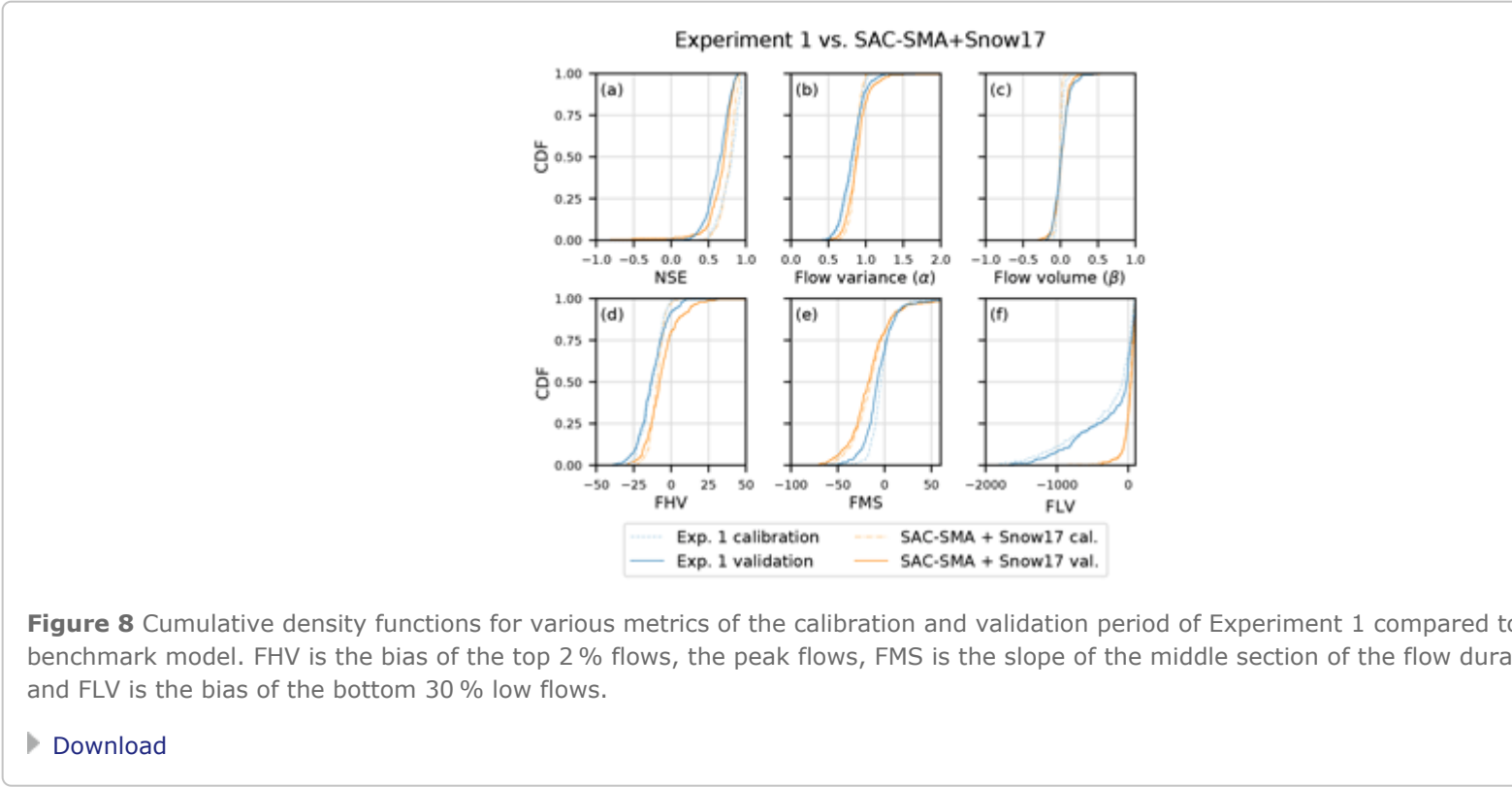
Figure 7a shows the spatial distribution of the LSTM performances for Experiment 1 in the validation period. In over 50 % of the catchments an NSE of 0.65 or above is found, with a mean NSE of 0.63 over all catchments. We can see that the LSTM performs better in catchments with snow influence (New England and Pacific Northwest regions) and catchments with higher mean annual precipitation (also the New England and Pacific Northwest regions, but also basins in the western part of the Arkansas-White-Red region; see Fig. 5a for precipitation distribution). The performance deteriorates in the more arid catchments, which are located in the western part of the Arkansas-White-Red region, where no discharge is observed for longer periods of the year (see Fig. 5b). Having a constant value of discharge (zero) for a high percentage of the training samples seems to be difficult information for the LSTM to learn and to reproduce this hydrological behaviour. However, if we compare the results for these basins to the benchmark model (Fig. 7b), we see that for most of these catchments the LSTM outperforms the latter, meaning that the benchmark model did not yield satisfactory results for these catchments either. In general, the visualization of the differences in the NSE shows that the LSTM performs slightly better in the northern, more snow-influenced catchments, while the SAC-SMA + Snow-17 performs better in the catchments in the south-east. This clearly shows the advantage of using LSTMs, since the snow accumulation and snowmelt processes are correctly reproduced, despite their inherent complexity. This suggests that the model learns these long-term dependencies, i.e. the time lag between precipitation falling as snow during the winter and runoff generation in spring with warmer temperatures. The median value of the NSE differences is  $-0.03$ , which means that the benchmark model slightly outperforms the LSTM. Based on the mean NSE value (0.58 for the benchmark model, compared to 0.63 for the LSTM of this Experiment), the LSTM outperforms the benchmark results.



**Figure 7** Panel **(a)** shows the NSE of the validation period of the models from Experiment 1 and panel **(b)** the difference of the NSE between the LSTM and the benchmark model (blue colours ( $>0$ ) indicate that the LSTM performs better than the benchmark model and ( $<0$ ) the other way around). The colour maps are limited to  $[0, 1]$  for the NSE and  $[-0.4, 0.4]$  for the NSE differences for better visualization.

► [Download](#)

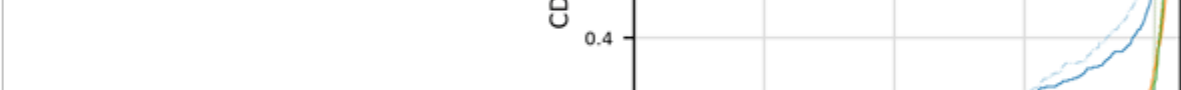
In Fig. 8, we present the cumulative density functions (CDF) for various metrics for the calibration and validation period. We see that the LSTM and the benchmark model work comparably well for all but the FLV (bias of the bottom 30 % low flows) metric. The underestimation of the peak flow in both models could be expected when using the MSE as the objective function for calibration (Gupta et al., 2009), as the LSTM underestimates the peaks more strongly compared to the benchmark model (Fig. 8d). In contrast, the middle section of the flow duration curve is better represented in the LSTM (Fig. 8e). Regarding the performance in terms of the NSE, the LSTM shows fewer negative outliers compared to the benchmark model, which seems to be more robust. The poorest model performance in the validation period is an NSE of  $-0.42$  compared to  $-20.68$  of the benchmark model (Snow-17). Figure 8f shows large differences between the LSTM and the SAC-SMA + Snow-17 model regarding the FLV metric. The FLV is sensitive to the one single minimum flow in the time series, since it compares the area between the FDC and this minimum value to the space of the observed and simulated discharge. The discharge from the LSTM model, which has no exponential outflow function like traditional hydrological models, can easily drop to diminutive numbers or even zero, to which we limited our model output. A rational solution for this issue is to introduce just one additional parameter and to limit the simulated discharge not to zero, but to the minimum observed flow from the calibration period. Figure 9 shows the effect of this approach on the CDF of the FLV. We can see that this solution leads to better FLV values compared to the benchmark model. Other metrics, such as the NSE, are almost unaffected by this change, since these low-flow values only marginally influence the resulting NSE values (not shown here).



**Figure 8** Cumulative density functions for various metrics of the calibration and validation period of Experiment 1 compared to the benchmark model. FHV is the bias of the top 2 % flows, the peak flows, FMS is the slope of the middle section of the flow duration curve, and FLV is the bias of the bottom 30 % low flows.

► [Download](#)

From the CDF of the NSE in Fig 8a, we can also observe a trend towards higher values in the calibration compared to the validation for both modelling approaches. This is a sign of overfitting, and in the case of the LSTM, could be tackled by a smaller network size, regularization or more data. However, we want to highlight again that achieving the best model performance possible was not the goal of this study, but rather testing the general ability of the LSTM to reproduce runoff processes.

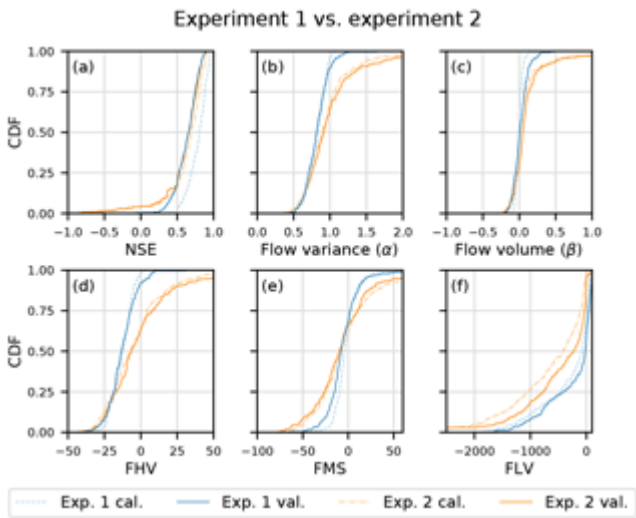


### 3.3 LSTMs as regional hydrological models

We now analyse the results of the four regional models that we trained for the four investigated HUCs in Experiment 2.

Figure 10 shows the difference in the NSE between the model outputs from Experiments 1 and 2. For some basins, the regional models perform significantly worse (dark red) than the individually trained models from Experiment 1. However, from the histograms of the differences we can see that the median is almost zero, meaning that in 50 % of the basins the regional model performs better than the specifically trained for a single basin. Especially in the New England region the regional model performed better for almost all basins (for two in the far north-east). In general, for all HUCs and catchments, the median difference is  $-0.001$ .

From Fig. 11 it is evident that the increased data size of the regional modelling approach (Experiment 2) helps to attenuate the performance difference between the calibration and validation periods, which could be observed in Experiment 1 probably as a result of overfitting. From the CDF of the NSE (Fig. 11a) we can see that Experiment 2 performed worse for approximately 20 % of the basins, while being slightly better for the remaining watersheds. We can also observe that the regional models show a more balanced under- and over-estimation, while the models from Experiment 1 as well as the benchmark model tend to underestimate the discharge (see Fig. 10). This is not too surprising, since we train one model on a specific basin with different discharge characteristics, where the model minimizes the error between simulated and observed discharge for all basins at the same time. On average, the regional model will therefore equally over- and under-estimate the observed discharge.



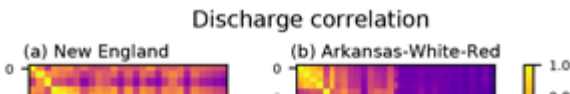
**Figure 11** Cumulative density functions for several metrics of the calibration and validation period of the models from Experiment 1 compared to the regional models from Experiment 2. FHV is the bias of the top 2 % flows, the peak flows, FMS is the slope of the middle section of the flow duration curve and FLV is the bias of the bottom 30 % low flows.

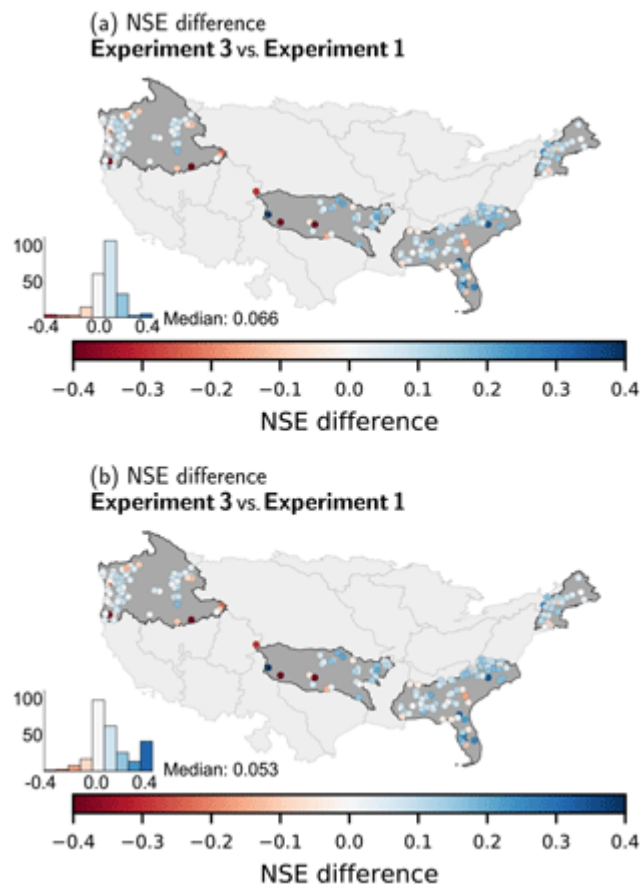
► [Download](#)

The comparison of the performances of Experiment 1 and 2 shows no clear consistent pattern for the investigated HUCs, but reveals a trend toward higher NSE values in the New England region and to lower NSE values in the Arkansas-White-Red region. The reason for these differences might become clearer once we look at the correlation in the observed discharge time series of the basins within both regions (see Fig. 12). We can see that in the New England region (where the regional model performed better for most of the catchments compared to the individual models of Experiment 1) many basins have a strong correlation in their discharge time series. Conversely, for the Arkansas-White-Red region the overall image of the correlation plot is much different. While some basins exist in the eastern part of the HUC with high discharge correlation, especially the basins in the western, more arid part have no inter-correlation at all. The results suggest that a single LSTM calibrated on a group of basins could generally be better in predicting the discharge of a group of basins compared to many LSTMs trained separately for each of the basins within the group especially when the group's basins exhibit a strong correlation in their discharge behaviour.

### 3.4 The effect of fine-tuning

In this section, we analyse the effect of fine-tuning the regional model for a few number of epochs to a specific catchment.



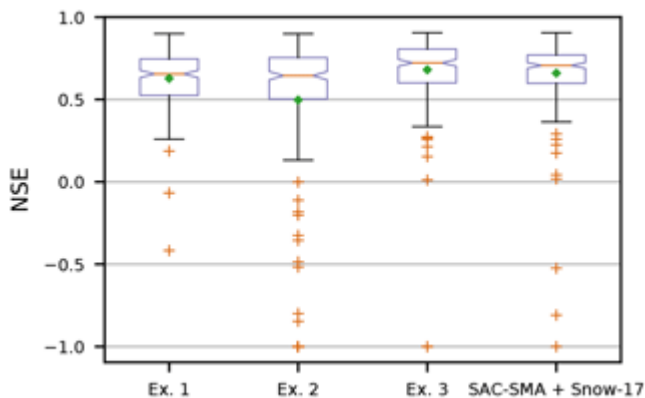


**Figure 13** Panel (a) shows the difference of the NSE in the validation period of Experiment 3 compared to the models of Experiment 1 and panel (b) in comparison to the models of Experiment 2. Blue colours ( $>0$ ) indicate in both cases that the fine-tuned model of Experiment 3 perform better and red colours ( $<0$ ) the opposite. The NSE differences are capped at  $[-0.4, 0.4]$  for better visualization.

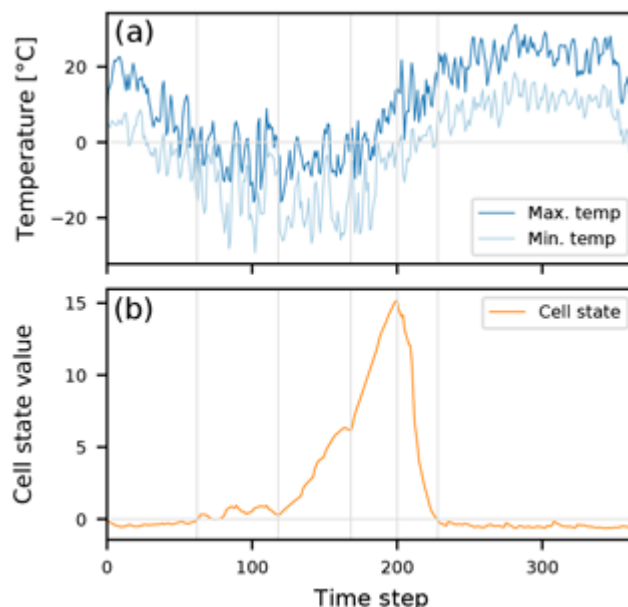
[Download](#)

### 3.5 A hydrological interpretation of the LSTM

To round off the discussion of this manuscript, we want to come back to the LSTM and try to explain it again in comparison to the structure of a classical hydrological model. Similar to continuous hydrological models, the LSTM processes the input data time step after time step. Every time step, the input data (here meteorological forcing data) are used to update a number of values in the LSTM internal cell states. In comparison to traditional hydrological models, the cell states can be interpreted as storages that are often used for e.g. snow accumulation, soil water content, or groundwater storage. Updating the internal cell states (or storages) is regulated through a number of so-called gates: one that regulates the depletion of the storages, a second that regulates the increase in the storages and a third that regulates the output from the storages. Each of these gates comes with a set of adjustable parameters that are adapted during a calibration period (referred to as *training*). During the validation period, updates of the cell states depend only on the input at a specific time step and the states from the previous time step (given the *learned* parameters of the calibration period).



**Figure 14** Boxplot of the NSE of the validation period for our three Experiments and the benchmark model. The NSE is capped at  $[-1, 1]$  for better visualization.



**Figure 15** Evolution of a specific cell state in the LSTM **(b)** compared to the daily min and max temperature, with accumulation and depletion in spring **(a)**. The vertical grey lines are included for better guidance.

► [Download](#)

## 4 Summary and conclusion

This contribution investigated the potential of using Long Short-Term Memory networks (LSTMs) for simulating runoff from meteorological observations. LSTMs are a special type of recurrent neural networks with an internal memory that has the ability to learn and store dependencies of the input–output relationship. Within three experiments, we explored possible applications of LSTMs and demonstrated that they are able to simulate the runoff with competitive performance compared to a baseline hydrological model (here the SAC-SMA model). In the first experiment we looked at classical single basin modelling, in a second experiment we trained one model for each of the regions we investigated, and in a third experiment we showed that using a pre-trained model helps to increase the model performance in single basins. Additionally, we showed an illustrative example why traditional RNNs should be avoided in favour of LSTMs for the task of predicting runoff from meteorological observations.

The goal of this study was to explore the potential of the method and not to obtain the best possible realization of the LSTM model for a catchment (see Sect. 2.5). It is therefore very likely that better performing LSTMs can be found by an exhaustive (catchment-wide) hyperparameter search. However, with our simple calibration approach, we were already able to obtain comparable (or even slightly better) model performances compared to the well-established SAC-SMA + Snow-17 model.

In summary, the major findings of the present study are the following.

- LSTMs are able to predict runoff from meteorological observations with accuracies comparable to the well-established SAC-SMA + Snow-17 model.
- The 15 years of daily data used for calibration seem to constitute a lower bound of data requirements.
- Pre-trained knowledge can be transferred into different catchments, which might be a possible approach for reducing the data demand and/or regionalization applications, as well as for prediction in ungauged basins or basins with few observations.

The data intensive nature of the LSTMs (as for any deep learning model) is a potential barrier for applying them in data-scarce regions (or for the usage within a single basin with limited data). We do believe that the use of “pre-trained LSTMs” (as explored in Experiment 3) is a promising way to reduce the large data demand for an individual basin. However, further research is needed to verify this hypothesis. Ultimately, however, LSTMs will always strongly rely on the available data for calibration. Thus, even if less data are needed, it could be a disadvantage in comparison to physically based models, which – at least in theory – are not reliant on calibration and can thus be applied with ease to new situations or catchments. However, more and more large-sample data sets are emerging which will catalyse further applications of LSTMs. In this context, it is also imaginable that adding physical catchment properties as an additional input layer to the LSTM may enhance the predictive power and ability of LSTMs to work as regional models and to make predictions in ungauged basins.

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., M., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vi, Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X.: TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, available at: <https://www.tensorflow.org/> (last access: 21 November 2018), 2016. [a](#)

Abraham, R. J., Anctil, F., Coulibaly, P., Dawson, C. W., Mount, N. J., See, L. M., Shamseldin, A. Y., Solomatine, D. P., Toth, E., and R. L.: Two decades of anarchy? Emerging themes and outstanding challenges for neural network river forecasting, *Prog. Phys. Geogr.*, 480–513, 2012. [a](#)

Adams, T. E. and Pagaon, T. C. (Eds.): *Flood Forecasting: A Global Perspective*, Academic Press, Boston, MA, USA, 2016. [a](#)

Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-scale studies, *Hydrol. Earth Syst. Sci.*, 21, 5293–5313, <https://doi.org/10.5194/hess-21-5293-2017>, 2017a. [a](#), [b](#), [c](#)

Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: Catchment attributes for large-sample studies, UCAR/NCAR, Boulder, CO, <https://doi.org/10.5065/D6G73C3Q>, 2017b. [a](#), [b](#)

Anderson, E. A.: National Weather Service River Forecast System - Snow Accumulation and Ablation Model, Tech. Rep. November 1973, Department of Commerce, Silver Spring, USA, 1973. [a](#)

ASCE Task Committee on Application of Artificial Neural Networks: Artificial Neural Networks in Hydrology. II: Hydrologic Applications, *J. Hydraul. Eng.*, 126, 124–137, 2000. [a](#)

Assem, H., Ghariba, S., Makrai, G., Johnston, P., Gill, L., and Pilla, F.: Urban Water Flow and Water Level Prediction Based on Deep Learning, in: *ECML PKDD 2017: Machine Learning and Knowledge Discovery in Databases*, 317–329, Springer, Cham., 2017. [a](#)

Bengio, Y.: Practical recommendations for gradient-based training of deep architectures, in: *Neural networks: Tricks of the trade*, Springer, Berlin, Heidelberg, 2012. [a](#)

Bengio, Y., Simard, P., and Frasconi, P.: Learning long-term dependencies with gradient descent is difficult, *IEEE T. Neural Networks*, 5, 166, 1994. [a](#), [b](#)

Beven, K.: Linking parameters across scales: subgrid parameterizations and scale dependent hydrological models, *Hydrol. Processes*, 9, 525, 1995. [a](#)

Beven, K.: *Rainfall-Runoff Modelling: The Primer*, John Wiley & Sons, Chichester, UK, 2001. [a](#)

Blöschl, G., Sivapalan, M., Wagener, T., Viglione, A., and Savenije, H. (Eds.): *Runoff Prediction in Ungauged Basins: Synthesis of Processes, Places and Scales*, Cambridge University Press, UK, 465 pp., 2013. [a](#)

Burnash, R. J. C., Ferral, R. L., and McGuire, R. A.: A generalised streamflow simulation system conceptual modelling for digital simulation, Tech. rep., US Department of Commerce National Weather Service and State of California Department of Water Resources, Sacramento, USA, 1973. [a](#)

Buytaert, W. and Beven, K.: Regionalization as a learning process, *Water Resour. Res.*, 45, 1–13, 2009. [a](#)

Carriere, P., Mohaghegh, S., and Gaskar, R.: Performance of a Virtual Runoff Hydrographic System, *Water Resources Planning and Management*, 122, 120–125, 1996. [a](#)

Chollet, F.: Keras, available at: <https://github.com/fchollet/keras> (last access: 1 April 2018), 2015. [a](#)

Clark, M. P., Bierkens, M. F. P., Samaniego, L., Woods, R. A., Uijlenhoet, R., Bennett, K. E., Pauwels, V. R. N., Cai, X., Wood, A., and Peters-Lidard, C. D.: The evolution of process-based hydrologic models: historical challenges and the collective quest for physical realism, *Hydrol. Earth Syst. Sci.*, 21, 3427–3440, <https://doi.org/10.5194/hess-21-3427-2017>, 2017. [a](#)

Daniell, T. M.: Neural networks. Applications in hydrology and water resources engineering, in: *Proceedings of the International Conference on Water Resource Symposium*, vol. 3, 797–802, Institution of Engineers, Perth, Australia, 1991. [a](#)

Duan, Q., Gupta, V. K., and Sorooshian, S.: Shuffled complex evolution approach for effective and efficient global minimization, *J. Geophys. Res.*, 98, 501–521, 1993. [a](#)

Fang, K., Shen, C., Kifer, D., and Yang, X.: Prolongation of SMAP to Spatiotemporally Seamless Coverage of Continental U.S. Using a Deep Learning Neural Network, *Geophys. Res. Lett.*, 44, 11030–11039, 2017. [a](#)

Farabet, C., Couprie, C., Najman, L., and Lecun, Y.: Learning Hierarchical Features for Scene Labeling, *IEEE T. Pattern Anal.*, 35, 1511–1527, 2013. [a](#)

Hestness, J., Narang, S., Ardalani, N., Damos, G. F., Jun, H., Kianinejad, H., Patwary, M. M. A., Yang, Y., and Zhou, Y.: Deep Learning for Rainfall–Runoff Modelling: A Review, *Water Resour. Res.*, 53, 1–15, 2017. <https://arxiv.org/abs/1712.00409> (last access: 21 November 2018), 2017. [a](#)

Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-R., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., and Kingsbury, C.: Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups, *IEEE Signal Process. Mag.*, 35, 82–97, 2012. [a](#)

Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, *Neural Comput.*, 9, 1735–1780, 1997. [a](#), [b](#), [c](#), [d](#)

Hsu, K., Gupta, H. V., and Soroochian, S.: Application of a recurrent neural network to rainfall-runoff modeling, *Proc., Aesthetics and Construction Environment, ASCE, New York*, 68–73, 1997. [a](#), [b](#)

Hunter, J. D.: Matplotlib: A 2D graphics environment, *Comput. Sci. Eng.*, 9, 90–95, 2007. [a](#)

Kirchner, J. W.: Getting the right answers for the right reasons: Linking measurements, analyses, and models to advance the science of hydrology, *Water Resour. Res.*, 42, 1–5, 2006. [a](#)

Kollet, S. J., Maxwell, R. M., Woodward, C. S., Smith, S., Vanderborght, J., Vereecken, H., and Simmer, C.: Proof of concept of regional hydrologic simulations at hydrologic resolution utilizing massively parallel computer resources, *Water Resour. Res.*, 46, 1–7, 2010. [a](#)

Krizhevsky, A., Sutskever, I., and Hinton, G. E.: ImageNet Classification with Deep Convolutional Neural Networks, *Adv. Neur. Inf. Process.*, 25, 1105, 2012. [a](#)

Kumar, D. N., Raju, K. S., and Sathish, T.: River Flow Forecasting using Recurrent Neural Networks, *Water Resour. Manag.*, 18, 1–10, 2004. [a](#)

LeCun, Y. A., Bottou, L., Orr, G. B., and Müller, K. R.: Efficient backprop, Springer, Berlin, Heidelberg, Germany, 2012. [a](#), [b](#)

Lindström, G., Pers, C., Rosberg, J., Strömquist, J., and Arheimer, B.: Development and testing of the HYPE (Hydrological Prediction Environment) water quality model for different spatial scales, *Hydrol. Res.*, 41, 295–319, 2010. [a](#)

Marçais, J. and de Dreuzay, J. R.: Prospective Interest of Deep Learning for Hydrological Inference, *Groundwater*, 55, 688–692, 2017. [a](#)

Maurer, E. P., Wood, A. W., Adam, J. C., Lettenmaier, D. P., and Nijssen, B.: A long-term hydrologically based dataset of land surface states and states for the conterminous United States, *J. Climate*, 15, 3237–3251, 2002. [a](#)

McKinney, W.: Data Structures for Statistical Computing in Python, Proceedings of the 9th Python in Science Conference, Location: Austin, Texas, USA, 1697900, 51–56, available at: <http://conference.scipy.org/proceedings/scipy2010/mckinney.html> (last access: 1 April 2010). [a](#)

Merz, R., Blöschl, G., and Parajka, J.: Regionalisation methods in rainfall-runoff modelling using large samples, *Large Sample Basin Analysis: Experiments for Hydrological Model Parameterization: Results of the Model Parameter Experiment–MOPEX*, IAHS Publ., 307, 117–124, 2005. [a](#)

Minns, A. W. and Hall, M. J.: Artificial neural networks as rainfall-runoff models, *Hydrolog. Sci. J.*, 41, 399–417, 1996. [a](#)

Mu, Q., Zhao, M., and Running, S. W.: Improvements to a MODIS global terrestrial evapotranspiration algorithm, *Remote Sens. Environ.*, 115, 1781–1800, 2011. [a](#)

Mulvaney, T. J.: On the use of self-registering rain and flood gauges in making observations of the relations of rainfall and of flood discharge in a given catchment, in: *Proceedings Institution of Civil Engineers, Dublin*, Vol. 4, 18–31, 1850. [a](#)

Myneni, R. B., Hoffman, S., Knyazikhin, Y., Privette, J. L., Glassy, J., Tian, Y., Wang, Y., Song, X., Zhang, Y., Smith, G. R., Los, S. O., Morisette, J. T., Votava, P., Nemani, R. R., and Running, S. W.: Global products of vegetation leaf area and fraction absorbed PAR from year one of MODIS data, *Remote Sens. Environ.*, 83, 214–231, 2002. [a](#)

Nash, J. E. and Sutcliffe, J. V.: River Flow Forecasting Through Conceptual Models Part I—a Discussion of Principles, *J. Hydrol.*, 10, 282–290, 1970. [a](#)

Newman, A., Sampson, K., Clark, M., Bock, A., Viger, R., and Blodgett, D.: A large-sample watershed-scale hydrometeorological dataset for the contiguous USA, UCAR/NCAR, Boulder, CO, USA, <https://doi.org/10.5065/D6MW2F4D>, 2014. [a](#), [b](#)

Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., Viger, R. J., Blodgett, D., Brekke, L., Arnold, J. R., Hoppe, J., and Duan, Q.: 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, *Hydrol. Earth Syst. Sci.*, 19, 209–223, <https://doi.org/10.5190/pd-2015-209-2015>, 2015. [a](#), [b](#), [c](#), [d](#)



Shen, C., Laloy, E., Elshorbagy, A., Albert, A., Bales, J., Chang, F.-J., Ganguly, S., Hsu, K.-L., Kifer, D., Fang, Z., Fang, K., Li, D., Tsai, W.-P.: HESS Opinions: Incubating deep-learning-powered hydrologic science advances as a community, *Hydrol. Earth Syst. Sci.*, 5639–5656, <https://doi.org/10.5194/hess-22-5639-2018>, 2018. [a](#)

Shi, X., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-K., and Woo, W.-C.: Convolutional LSTM network: A machine learning approach for precipitation nowcasting, *Adv. Neur. In.*, 28, 802–810, 2015. [a](#)

Sivapalan, M.: Prediction in ungauged basins: a grand challenge for theoretical hydrology, *Hydrol. Process.*, 17, 3163–3170, 2002. [a](#)

Solomatine, D., See, L. M., and Abrahart, R. J.: Data-driven modelling: concepts, approaches and experiences, in: *Practical hydroinformatics*, 17–30, Springer, Berlin, Heidelberg, 2009. [a](#)

Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, *J. Mach. Learn. Res.*, 15, 1929–1958, 2014. [a](#)

Stanzel, P., Kahl, B., Haberl, U., Herrnegger, M., and Nachtnebel, H. P.: Continuous hydrological modelling in the context of real-time forecasting in alpine Danube tributary catchments, *IOP C. Ser. Earth Env.*, 4, 012005, <https://doi.org/10.1088/1755-1307/4/1/012005>, 2008. [a](#)

Sutskever, I., Vinyals, O., and Le, Q. V.: Sequence to sequence learning with neural networks, in: *Advances in neural information processing systems*, 3104–3112, 2014. [a](#)

Tao, Y., Gao, X., Hsu, K., Sorooshian, S., and Ihler, A.: A Deep Neural Network Modeling Framework to Reduce Bias in Satellite Rainfall Products, *J. Hydrometeorol.*, 17, 931–945, 2016. [a](#)

Thielen, J., Bartholmes, J., Ramos, M.-H., and de Roo, A.: The European Flood Alert System – Part 1: Concept and development, *Earth Syst. Sci.*, 13, 125–140, <https://doi.org/10.5194/hess-13-125-2009>, 2009. [a](#)

Thornton, P. E., Thornton, M. M., Mayer, B. W., Wilhelmi, N., Wei, Y., Devarakonda, R., and Cook, R.: Daymet: Daily surface weather and meteorological data, 1 km grid for North America, 1980–2008, Oak Ridge National Laboratory (ORNL) Distributed Active Archive Center for Biogeochemistry (DAAC), Oak Ridge, Tennessee, USA, 2012. [a](#)

Tompson, J., Jain, A., LeCun, Y., and Bregler, C.: Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation, in: *Proceedings of Advances in Neural Information Processing Systems*, 27, 1799–1807, 2014. [a](#)

Van Der Walt, S., Colbert, S. C., and Varoquaux, G.: The NumPy array: A structure for efficient numerical computation, *Comput. Sci.*, 13, 22–30, 2011. [a](#)

van Rossum, G.: Python tutorial, Technical Report CS-R9526, Tech. rep., Centrum voor Wiskunde en Informatica (CWI), Amsterdam, Netherlands, 1995. [a](#)

Wesemann, J., Herrnegger, M., and Schulz, K.: Hydrological modelling in the anthroposphere: predicting local runoff in a heavily forested high-alpine catchment, *J. Mt. Sci.*, 15, 921–938, 2018. [a](#)

Wood, E. F., Roundy, J. K., Troy, T. J., van Beek, L. P. H., Bierkens, M. F. P., Blyth, E., de Roo, A., Döll, P., Ek, M., Famiglietti, J., Gochis, D. J., van de Giesen, N., Houser, P., Jaffé, P. R., Kollet, S., Lehner, B., Lettenmaier, D. P., Peters-Lidard, C., Sivapalan, M., Sheffield, J., and Whitehead, P.: Hyperresolution global land surface modeling: Meeting a grand challenge for monitoring Earth's terrestrial water resources, *Resour. Res.*, 47, W05301, <https://doi.org/10.1029/2010WR010090>, 2011. [a](#)

Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H., Meng, J., and Livneh, B.: Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2), 1. Intercomparison and application of model products, *J. Geophys. Res.-Atmos.*, 117, D03109, <https://doi.org/10.1029/2011JD001617>, 2012. [a](#)

Yilmaz, K. K., Gupta, H. V., and Wagener, T.: A process-based diagnostic approach to model evaluation: Application to the NWS National Water Research Institute hydrologic model, *Water Resour. Res.*, 44, 1–18, 2008. [a](#)

Yosinski, J., Clune, J., Bengio, Y., and Lipson, H.: How transferable are features in deep neural networks?, *Adv. Neur. In.*, 27, 1–9, 2015. [a](#)

Young, P. C. and Beven, K. J.: Data-based mechanistic modelling and the rainfall-flow non-linearity, *Environmetrics*, 5, 335–363, 2003. [a](#)

Zhang, D., Lindholm, G., and Ratnaweera, H.: Use long short-term memory to enhance Internet of Things for combined sewer overflow monitoring, *J. Hydrol.*, 556, 409–418, 2018. [a](#)

Zhang, J., Zhu, Y., Zhang, X., Ye, M., and Yang, J.: Developing a Long Short-Term Memory (LSTM) based model for predicting water depth in agricultural areas, *J. Hydrol.*, 561, 918–929, 2018. [a](#)