

Journal Pre-proofs

Research papers

Improving LSTM hydrological modeling with spatiotemporal deep learning and multi-task learning: a case study of three mountainous areas on the Tibetan Plateau

Bu Li, Ruidong Li, Ting Sun, Aofan Gong, Fuqiang Tian, Mohd Yawar Ali Khan, Guangheng Ni

PII: S0022-1694(23)00343-8
DOI: <https://doi.org/10.1016/j.jhydrol.2023.129401>
Reference: HYDROL 129401

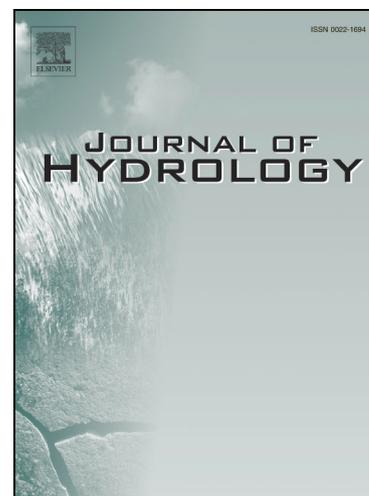
To appear in: *Journal of Hydrology*

Received Date: 13 April 2022
Revised Date: 10 March 2023
Accepted Date: 11 March 2023

Please cite this article as: Li, B., Li, R., Sun, T., Gong, A., Tian, F., Yawar Ali Khan, M., Ni, G., Improving LSTM hydrological modeling with spatiotemporal deep learning and multi-task learning: a case study of three mountainous areas on the Tibetan Plateau, *Journal of Hydrology* (2023), doi: <https://doi.org/10.1016/j.jhydrol.2023.129401>

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

© 2023 Published by Elsevier B.V.



1 **Improving LSTM hydrological modeling with**
2 **spatiotemporal deep learning and multi-task**
3 **learning: a case study of three mountainous areas on**
4 **the Tibetan Plateau**

5

6 Bu Li¹, Ruidong Li¹, Ting Sun², Aofan Gong¹, Fuqiang Tian¹, Mohd Yawar Ali
7 Khan³, Guangheng Ni¹

8 ¹State Key Laboratory of Hydro-science and Engineering, Department of Hydraulic
9 Engineering, Tsinghua University, Beijing 100084, China

10 ²Institute for Risk and Disaster Reduction, University College London, London
11 WC1E 6BT, UK

12 ³Department of Hydrogeology, Faculty of Earth Sciences, King Abdulaziz University,
13 Jeddah 21589, Saudi Arabia

14 Correspondence to: G. Ni, ghni@tsinghua.edu.cn

15

16 **Abstract**

17 Long short-term memory (LSTM) networks have demonstrated their excellent
18 capability in processing long-length temporal dynamics and have proven to be
19 effective in precipitation-runoff modeling. However, the current LSTM hydrological
20 models lack the incorporation of multi-task learning and spatial information, which
21 limits their ability to make full use of meteorological and hydrological data. To
22 address this issue, this study proposes a spatiotemporal deep-learning (DL)-based
23 hydrological model that couples the 2-Dimension convolutional neural network
24 (CNN) and LSTM and introduces actual evaporation (E_a) as an additional training
25 target. The proposed CNN-LSTM model is tested on three large mountainous basins
26 on the Tibetan Plateau, and the results are compared to those obtained from the
27 LSTM-only model. Additionally, a probe method is used to decipher the internal
28 embedding layers of the proposed DL models. The results indicate that both LSTM
29 and CNN-LSTM hydrological models perform well in simulating runoff (Q) and E_a ,
30 with Nash-Sutcliffe efficiency coefficients ($NSEs$) higher than 0.82 and 0.95,

31 respectively. The higher $NSEs$ suggest that introducing spatial information into
32 LSTM-only models can improve the overall and peak model performance. Moreover,
33 multi-task simulation with LSTM-only models shows better accuracy in the
34 estimation of Q volume and performance, with $NSEs$ increasing by approximately
35 0.02. The probe method also reveals that CNN can capture the basin-averaged
36 meteorological values in CNN-LSTM models, while LSTM Q (E_a) models contain
37 the information about the known E_a (Q) process. Overall, this study demonstrates the
38 value of spatial information and multi-task learning in LSTM hydrological modeling
39 and provides a perspective for interpreting the internal embedding layers of DL
40 models.

41

42 **Highlight**

- 43 (1) Spatiotemporal DL model enhances LSTM by introducing spatial information.
- 44 (2) Multi-task simulation improves LSTM with enhanced performance in estimating
45 Q volume.
- 46 (3) Probe method shows CNN captures basin-averaged meteorological values in
47 CNN-LSTM, LSTM Q (E_a) models contain known process.

48 **Keywords**

49 CNN-LSTM; spatiotemporal; multi-task; actual evaporation; Tibetan Plateau

50

51 1. Introduction

52 Hydrological models, either physically-based or data-driven, play vital roles in flood
53 and drought disaster prevention, as well as in water resources management (Blöschl et
54 al. 2019). Recent advances in remote sensing and computational techniques have
55 improved physically-based hydrological models by enhancing their capacity in
56 characterizing hydrodynamic processes at finer scales (Khatakho et al. 2021, Sood
57 and Smakhtin 2015, Tarek et al. 2020): notably, distributed hydrological models
58 (DHMs) can incorporate processes at a sub-basin or other calculation unit scale (Li et
59 al. 2021a). However, the application of DHMs is still limited owing to the unsolved
60 issue in scale mismatch (Blöschl et al. 2019, Blöschl and Sivapalan 1995, Gupta et al.
61 2008, Hrachowitz et al. 2013, Nearing et al. 2021).

62 Data-driven models can directly depict the statistical relationship between inputs and
63 outputs without explicit characterization of physical processes (Nearing et al. 2021).
64 Since the 1990s, machine learning (ML) techniques have been widely adopted in
65 hydrological modeling and have proven similar or better performance compared with
66 DHMs (Demirel et al. 2009, Hsu et al. 1995, Kratzert et al. 2018, Lees et al. 2021,
67 Nearing et al. 2021, Yang et al. 2020). In particular, deep learning (DL) models,
68 featuring neural networks, are the most commonly used ML techniques for
69 hydrological modeling (Nearing et al. 2021). They have evolved from original multi-
70 layer perceptron, i.e., artificial neural networks (ANNs, Yang et al. 2020), to more
71 advanced forms with enriched taxonomy, such as convolutional neural networks
72 (CNNs, Jiang et al. 2020), recurrent neural networks (RNNs, Sadeghi Tabas and
73 Samadi 2022), and its variant, long short-term memory networks (LSTM, Lees et al.
74 2021). ANNs, CNNs, and RNNs are among the earliest DL models with promising
75 performance for hydrological modeling (Demirel et al. 2009, Hsu et al. 1995, Khan et
76 al. 2019, Yang et al. 2020). As feed-forward neural networks, ANNs and CNNs
77 cannot directly represent temporal dynamics and are thus less able to accurately
78 characterize hydrological processes (Kratzert et al. 2018). In contrast, RNNs, which
79 process the input data chronologically by design, can consider temporal dynamics.
80 Although RNNs are more ideal for time series analysis, vanilla RNNs can hardly store
81 sequences over 10 time steps (Bengio et al. 1994) which limits their applicability in
82 modeling slow hydrological processes occurring at larger time scale, such as those
83 related to groundwater, snow, and glacier storage (Kratzert et al. 2018). The recently
84 emerging LSTM models can conquer such weakness with the unique internal gate
85 architectures and have demonstrated superior performance than the vanilla RNNs in
86 time series analysis (Hochreiter and Schmidhuber 1997). Since the first application in
87 hydrological modeling by Kratzert et al. (2018) for a precipitation-runoff (P - Q)
88 simulation in 530 American basins ($<2000 \text{ km}^2$), the LSTM P - Q models have been
89 used worldwide (e.g., 669 basins in Great Britain in Lees et al. 2021, Hanjiang River
90 in China in Liu et al. 2021) and have become one of the most powerful tools in P - Q

91 simulations.

92 Being promising in hydrological modeling, current DL-based hydrological models
93 still need improvements in three notable aspects (Nearing et al. 2021):

- 94 (1) Ability to resolve spatiotemporal features: hydrological processes modeling
95 depends heavily on the spatial patterns of meteorological forcing and
96 underlying surface characteristics. Yang et al. (2020) employed computer
97 vision to resolve spatial features in ANN P - Q modeling and demonstrated that
98 spatial information plays an important role in enhancing model robustness.
99 However, most existing studies about LSTM hydrological modeling almost
100 utilized basin spatially-averaged meteorological data as model inputs (e.g.,
101 Jiang et al. 2022, Kratzert et al. 2018, Lees et al. 2021), without fully
102 representing spatial features of inputs for the LSTM hydrological modeling.
103 Coupling 2-D CNN and LSTM is expected to bridge such gap by
104 simultaneously considering both temporal dynamics and spatial features (Miao
105 et al. 2019, Shi et al. 2015) and CNN-LSTM has proven to be promising in
106 different fields (e.g., P nowcasting (Miao et al. 2019, Shi et al. 2015) and
107 water quality forecast (Barzegar et al. 2020, Yang et al. 2021)).
- 108 (2) Consideration of multiple hydrological processes: differing from physically-
109 based hydrological models, LSTM-based models simulate the individual
110 hydrological process, such as Q (Feng et al. 2020, Frame et al. 2021),
111 groundwater (Ali et al. 2022, Nourani et al. 2022), and snow water equivalent
112 (Duan and Ullrich 2021) in most studies, but rarely simulate multiple
113 processes simultaneously. It makes LSTM-based hydrological models difficult
114 to explicitly consider the interactions between different hydrological processes
115 and diagnose models based on hydrological theories, such as the water balance
116 equation (Reichstein et al. 2019). Besides, some studies found that introducing
117 additional hydrological processes, such as the actual evaporation (denoted by
118 E_a in this work) process, in calibration for physical-based hydrological
119 models can enhance the Q simulation performance (Herman et al. 2018,
120 Nesru et al. 2020). Therefore, it is beneficial to investigate if considering
121 multiple hydrological processes simultaneously for LSTM-based hydrological
122 models can enhance their capacity in depicting more hydrological processes
123 and thus provide a more comprehensive diagnosis of hydrological variables.
- 124 (3) Physical interpretability: due to the “black-box” nature, DL-based
125 hydrological models have no explicit representation of physical processes and
126 thus remain being questioned by some hydrologists (Nearing et al. 2021). To
127 enhance the confidence of users and policymakers in adopting DL-based
128 hydrological models, improvements in the understanding of their physical

129 interpretability have been attempted recently (Arrieta et al. 2020, Jiang et al.
130 2022). For example, LSTM hydrological models have proven to learn a
131 generalizable representation of the underlying physical processes. LSTM
132 regional hydrological models outperform DHMs calibrated regionally, and
133 even calibrated for each basin individually (Feng et al. 2020, Kratzert et al.
134 2019, Sun et al. 2021). Besides, LSTM hydrological models are found to be
135 able to store the hidden information consistent with hydrological knowledge
136 (Jiang et al. 2022, Kratzert et al. 2018, Lees et al. 2022). However, the
137 physical interpretability of DL-models with respect to E_a process—a critical
138 component in the hydrological cycle—is yet to be investigated. Also, the
139 physical concepts of CNN outputs in CNN-LSTM models are still unclear.

140 In this work, we aim to overcome these deficiencies by developing a spatiotemporal
141 DL-based hydrological model by coupling 2-D CNN and LSTM (CNN-LSTM) and
142 introducing multi-task learning. The potential of simultaneous multi-task (MT)
143 learning in DL-based hydrological models is also investigated by involving E_a
144 process as additional learning target. Besides, we also advance the understanding of
145 physical interpretability of DL-based models by extracting the meteorological and
146 hydrological processes hidden in the proposed LSTM and CNN-LSTM models.

147 In the remainder of this paper, we first describe the proposed DL-based model by
148 introducing the basic architecture of CNN and LSTM and physical interpretability
149 method (Sec. 2), then evaluate the model performance in three large mountainous
150 basins on the Tibetan Plateau with comparison to the LSTM hydrological models
151 (Sec. 3); we also explore physical interpretations of this model with respect to the
152 hydrological and meteorological processes (Sec. 4).

153 2. Methods

154 2.1 Model development

155 We propose a DL-based hydrological model (Figure 1a) by coupling 2-D CNN
156 (Figure 1b) and LSTM (Figure 1c) to utilize their respective advantages: the former
157 for hierarchical spatial feature extraction while the latter for learning long temporal
158 dependencies. The proposed model can use 2-D spatial meteorological and underlying
159 surface data as input and predicts hydrological processes with daily Q and E_a as
160 output. We note, however, only meteorological data—daily P and mean temperature—
161 are used in this work as inclusion of surface characteristics demonstrate minimal
162 improvement in model performance (not shown). The model can perform either
163 single-task or multi-task learning by setting training targets: the former simulates the
164 individual hydrological process, while the latter focuses on two or more processes
165 simultaneously. Below we focus on the model design and description of key
166 components; more technical details refer to Appendix A.

167 **2.1.1 Convolutional neural networks (CNNs)**

168 CNNs (Figure 1b), are a particular type of feed-forward neural network, including the
169 input, convolutional, pooling, and full connection layer (LeCun et al. 1998). The
170 convolutional layer, the core of CNNs, uses convolutional kernels to extract
171 information from various N-Dimensions model inputs. We utilize 2-D CNN to
172 capture the spatial information of meteorological data. The convolutional layer
173 processes meteorological data by reducing the spatial size (width and height) of inputs
174 (features), increasing the channel number, and generating the 1-dimension sequence
175 finally. Figure 1b takes a study basin as the example to illustrate the data dimensions
176 in internal layers of the CNNs in this study and more details refer to Appendix A.

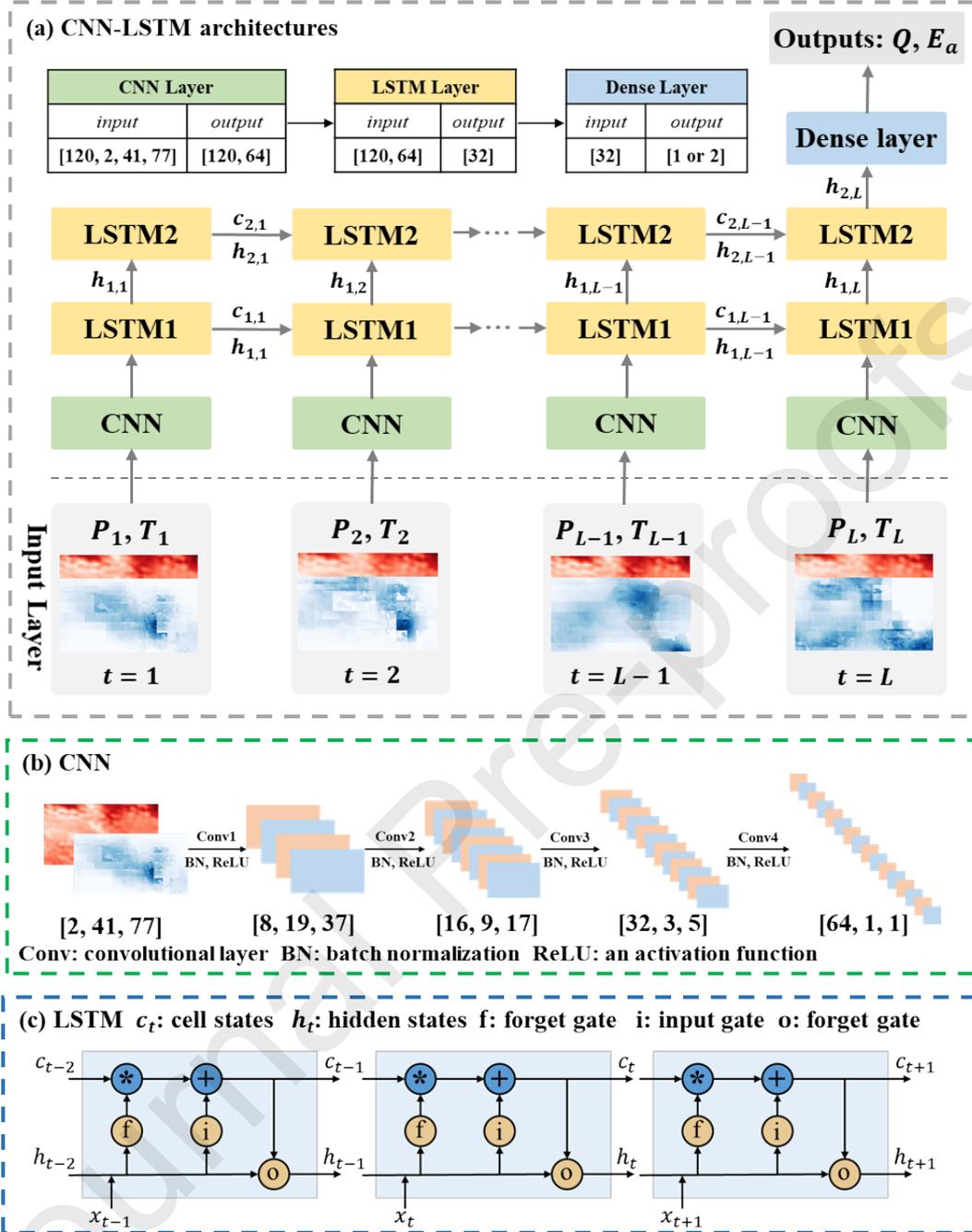


Figure 1. (a) CNN-LSTM model architectures and the dimensions of data before and after each layer (take the Yellow as an example and other basins are shown in Appendix A, same with Panel (b)) in this study. (b) The workflow of the CNN model and the dimensions of data in the internal layers of CNN model. (c) The internals of LSTM cells. The meaning of each unexplained variable is stated in the text.

177 **2.1.2 Long short-term memory (LSTM) networks**

178 The LSTM models (Figure 1c) are designed to alleviate the weakness of the vanilla
179 RNNs in processing long-length temporal dynamics (Hochreiter and Schmidhuber
180 1997, Sherstinsky 2020). The success of LSTM lies in the coordination of memory
181 cells (cell states; c_t in Figure 1c) and hidden cells (hidden states; h_t in Figure 1c) in
182 the internal architecture that capture the slow and quick evolution processes,
183 respectively. Besides, three gates (i.e., input, forget, and output) are designed to
184 control the information in each cell to be stored, removed, and passed, respectively.
185 These architectures are beneficial for LSTM models to process long-length temporal
186 dynamics. More detailed descriptions of LSTM within the context of hydrological
187 modeling refer to Kratzert et al. (2018).

188 **2.1.3 Multi-task learning**

189 Multi-task (MT) learning is to set multiple tasks as optimization targets in a DL-based
190 model (Caruana 1997). The training can benefit from the enriched representations of
191 MT information, thus enhancing the performance of each task using the information
192 enveloped in the other tasks. Besides, MT learning can achieve higher efficiency and
193 less over-fitting than single-task (ST) learning because it can lead the model to a more
194 general feature representation preferred by multiple related tasks (Li et al. 2022,
195 2023). The MT learning is introduced to investigate the effect of multi-hydrological-
196 process learning in LSTM-based hydrological modeling in this study.

197 **2.1.4 Coupling between CNN and LSTM**

198 To utilize the respective advantage of CNN and LSTM models, we develop a
199 spatiotemporal DL-based hydrological model by coupling CNN and LSTM models
200 (Figure 1a). At each time step, the CNN-LSTM coupling is realized via two stages:

- 201 1) the CNN reduces the 2-D gridded meteorological input into a 1-D sequence
202 and feed it to the LSTM1 layer;
- 203 2) the hidden states of the LSTM1 layer are input to the LSTM2 layer and cell
204 and hidden states of two LSTM layers are then fed to the corresponding layers
205 at the next time step, respectively.

206 After all processing through the CNN-LSTM coupled layer prior to time step L , the
207 outputs of the LSTM2 layer at the last time step are passed to the dense layer to obtain
208 the predictions: learned hydrological variables (Figure 1a). Depending on the target
209 mode—ST or MT—the DL model may produce different number of output variables.
210 ST model sets one hydrological process as optimization target while MT model sets
211 two or more hydrological processes as optimization targets by MT learning.

212 2.2 Physical interpretability method

213 The internal embedding layers of the DL hydrological models contain a large amount
214 of data that is not explicitly interpretable. These data may conceal some untrained
215 internal hydrological variables. For example, the internal embedding layers of Q
216 models may contain the information about E_a process. This study utilizes probes–
217 regression models that map the internal embedding layers of trained models to
218 untrained hydrological variables (Hewitt and Liang 2019, Lees et al. 2022)–to test
219 whether trained DL models can learn the known but untrained hydrological variables
220 and examine the internal representation and further physical interpretability of
221 models. The simplest form of probes is a linear regression (LR) model that connects
222 the learned embedding layers to a given output. This study explores untrained
223 hydrological variables from proposed LSTM and CNN-LSTM models based on LR
224 models and detailed experiment design refers to Sec. 4.1.

225 3. Model evaluation

226 We evaluate the performance of proposed CNN-LSTM model in three large
227 mountainous basins on the Tibetan Plateau (TP; Sec. 3.1). To assess the applicability
228 of CNN-LSTM in resolving spatial information in hydrological modeling, an LSTM-
229 only model is also configured as a benchmarking baseline: differing from the CNN-
230 LSTM model, the basin spatially-averaged meteorological data at each time step are
231 directly input to the LSTM-only model without the involvement of CNN model
232 (detailed in Kratzert et al. 2018). Also, ST and MT experiments (Sec. 3.2) are
233 designed to evaluate the MT learning performance of CNN-LSTM and LSTM models
234 and the influence of MT on DL-based hydrological modeling. In this study, the model
235 inputs are daily total P and mean air temperature (T) in line with similar studies
236 (e.g., Jiang et al. 2022). In addition to Q , we also select the E_a –an important
237 hydrological processes for which data are available– as the training target to evaluate
238 the effect of multi-task learning in LSTM modeling.

239 3.1 Study area and data

240 3.1.1 Study area

241 The Tibetan Plateau (TP; Figure 2a) is the highest plateau in the world, known as the
242 “Roof of the World”, or the “Third Pole”. Given the high altitude and vast glaciers of
243 TP, along with the many mighty rivers (e.g., Yangtze, Lancang/Mekong, Yarlung
244 Zangbo/Brahmaputra, and Yellow, among others) that provide enormous water
245 sources for downstream livelihoods and agricultural irrigation (Huss et al. 2017,
246 Immerzeel et al. 2010, Nan et al. 2021, Schaner et al. 2012, Wang et al. 2021, Zhang
247 et al. 2013), TP is also considered the “Water tower of Asia”. However, owing to the
248 incompleteness of knowledge of complex alpine hydrological processes, it is

249 challenging to adequately model hydrological processes by DHMs in the TP (Li et al.
250 2019b, Nan et al. 2021).

251 **Table 1.** Basic facts of the three study basins. “DEM” represents Digital Elevation
252 Model.

Basins	Area (km ²)	DEM range (m)	Average annual <i>P</i> (mm)	Average annual <i>Q</i> (10 ⁹ m ³)
Yellow	123,000	2,656-6,253	510	20
Yangtze	139,000	3,516-6,575	460	16
Lancang	91,000	1,243-6,334	830	32

253 To systematically evaluate the performance of the LSTM and CNN-LSTM
254 hydrological model in such a challenging environment, we select source regions of
255 three rivers (Figure 2a)–the Yellow River (Figure 2b), the Yangtze River (Figure 2c),
256 and the Lancang River (Figure 2d)–characterized as mountainous areas as the study
257 basins. The landforms of all three basins undulate greatly with elevation variability
258 greater than 3,000 m. The annual average *P* of three basins ranges from 460 to 830
259 mm. The detailed description of the three study basins is shown in Table 1. Yellow,
260 Yangtze, and Lancang are used hereinafter to denote their respective source regions if
261 not specified otherwise.

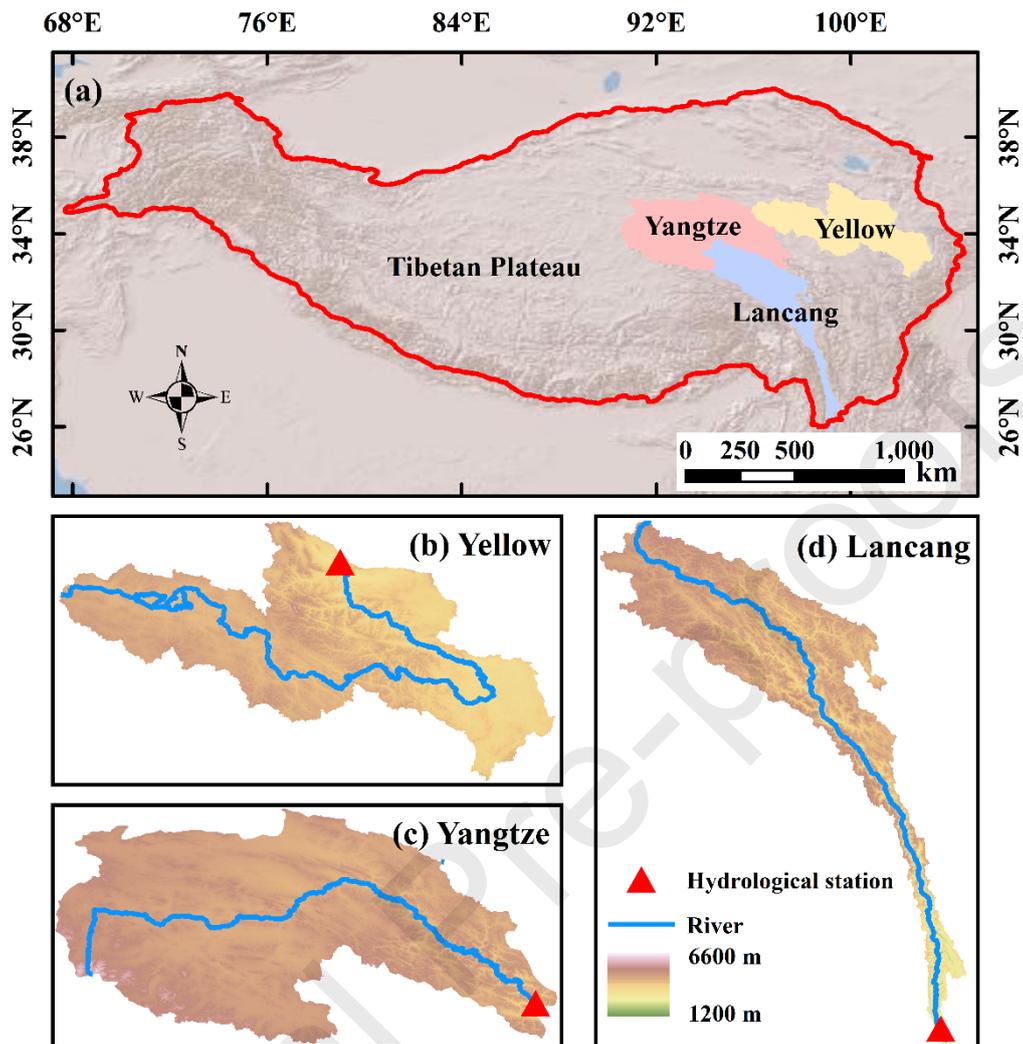


Figure 2. (a) The terrain of the Tibetan Plateau and the location of three study basins. The terrain and CNN-LSTM models input range of the Yellow River source region (b), the Yangtze River source region (c), and the Lancang River source region (d).

262 3.1.2 Data

263 Due to the limited availability of in situ observations in the study area, we utilize
 264 either the optimal remote sensing or reanalysis datasets for different input variables as
 265 follows:

- 266 (1) Precipitation (P): Multi-Source Weighted-Ensemble Precipitation (MSWEP)
 267 V2.2 with 0.1° spatial and 3h temporal resolution (Beck et al. 2017, Beck et al.
 268 2019).

269 (2) Air temperature (T): The air temperature at 2m AGL (T2) from the fifth
 270 generation of ECMWF atmospheric reanalysis of the global climate (ERA5)
 271 reanalysis dataset with 0.1° spatially and 1h temporal resolution (Hersbach et
 272 al. 2020)

273 **Table 2.** The length of the training, evaluation, and testing periods in three study
 274 basins.

Basins	Training period	Evaluation period	Testing period
Yellow	1983-2004	2006-2009	2011-2014
Yangtze	1983-2004	2006-2009	2011-2014
Lancang	1988-2004	2005-2007	2008-2010

275 While for output/evaluation dataset, we use the following datasets:

276 (1) Evaporation (E_a): Global Land Evaporation Amsterdam Model (GLEAM)
 277 v3.5a E_a dataset with 0.25° spatial and 1-day temporal resolution (Martens
 278 et al. 2017, Miralles et al. 2011). The basin spatially-averaged daily E_a are
 279 calculated as the model learning target.

280 (2) Runoff (Q): the daily in situ measurements collected at three hydrological
 281 stations (Figure 1) provided by local water agencies.

282 All the above datasets are pre-processed into daily resolution in line with the model
 283 time step and spilt into three periods for training, evaluation, and testing (Table 2).

284 3.2 Model evaluation experiment

285 3.2.1 Experiment design

286 To test the performance of LSTM and CNN-LSTM models for single-task (ST) and
 287 multi-task (MT) modes, a total of 18 scenarios are set up: 3 basins (Yellow, Yangtze,
 288 Lancang) \times 2 models (LSTM and CNN-LSTM) \times 3 tasks (Q , E_a , $Q + E_a$). Task Q
 289 and E_a are for the ST experiment and $Q + E_a$ for MT.

290 3.2.2 Evaluation metrics

291 This study utilizes Nash–Sutcliffe efficiency (NSE; Equation 1; J.E.Nash and

292 J.V.Sutcliffe 1970) and its three decompositions (Gupta et al. 2009), namely, the
 293 correlation coefficient (r ; Eq. 2), the variance bias (α ; Eq. 3), and the total volume
 294 bias (β ; Eq. 4) to evaluation metrics:

$$NSE = 1 - \frac{\sum_{t=1}^T (V_{sim} - V_{obs})^2}{\sum_{t=1}^T (V_{obs} - \overline{V_{obs}})^2} \quad (1)$$

$$r = \frac{\sum_{t=1}^T (V_{sim} - \overline{V_{sim}})(V_{obs} - \overline{V_{obs}})}{\sqrt{\sum_{t=1}^T (V_{sim} - \overline{V_{sim}})^2 \sum_{t=1}^T (V_{obs} - \overline{V_{obs}})^2}} \quad (2)$$

$$\alpha = \frac{\sigma_{sim}}{\sigma_{obs}} \quad (3)$$

$$\beta = \frac{(\mu_{sim} - \mu_{obs})}{\sigma_{obs}} \quad (4)$$

295 where V_{sim} (V_{obs}), σ_{sim} (σ_{obs}), and μ_{sim} (μ_{obs}) are the simulated (observed) values,
 296 standard deviations, and means, respectively. The variance (α) and total volume
 297 bias (β) measure the error in the standard deviation and the average values, respectively.

298 Besides, we calculate the peak, middle, and low values bias of the results grouped by
 299 the exceedance probability of observations to assess the range-specific model
 300 performance (Yilmaz et al. 2008):

$$BIV = \frac{\sum_{i=1}^I (V_{sim,i} - V_{obs,i})}{\sum_{i=1}^I V_{obs,i}} \quad (5)$$

301 where I is one of P , M , and L that represent the peak (0–0.02 exceedance
 302 probabilities), middle (0.3–0.7 exceedance probabilities), and low (0.7–1.0 exceedance
 303 probabilities) values, respectively.

304 3.3 Evaluation results

305 3.3.1 Performance of CNN-LSTM models

306 In general, CNN-LSTM Q models work remarkably well in all three study basins

307 with $NSEs > 0.89$ (Table 3) (Moriassi et al. 2007) and can accurately capture peak of
 308 Q with appropriate magnitudes and timing (Figure 3). The results are compared to
 309 some traditional hydrological models in other relevant studies, and it is indicated that
 310 CNN-LSTM Q models outperform them (see detailed in Appendix B). Besides,
 311 simulated E_a processes agree well with GLEAM data with $NSEs$ of 0.97 in all
 312 three study basins (Table 4; Figure 4-5). These results demonstrate that the proposed
 313 CNN-LSTM model performs favorably and thus is a reasonable approach to simulate
 314 hydrological processes in large mountainous basins.

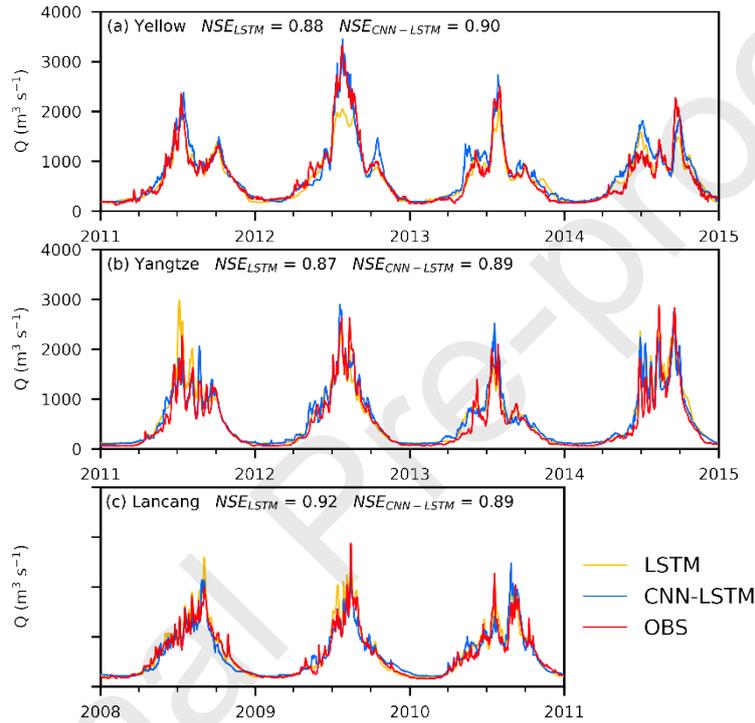


Figure 3. The comparison of simulated (ST Q models) and observed Q processes in the test period at three hydrological stations.

315 **Table 3.** The evaluation metrics of LSTM and CNN-LSTM hydrological models for
 316 Q simulation. LSTM (MT) and CNN-LSTM (MT) represent the LSTM and CNN-
 317 LSTM Q (multi-task) models, respectively.

NSE	r	α	β	BPV	BMV	BLV
-------	-----	----------	---------	-------	-------	-------

Yellow	LSTM	0.88	0.95	0.80	-0.06	-0.31	0.02	0.12
	CNN-LSTM	0.90	0.95	1.01	0.11	-0.00	0.17	0.16
	LSTM MT	0.90	0.95	0.85	-0.04	-0.29	0.03	0.10
	CNN-LSTM MT	0.86	0.94	0.98	0.15	-0.13	0.22	0.32
Yangtze	LSTM	0.87	0.94	1.01	0.12	-0.15	0.24	0.58
	CNN-LSTM	0.89	0.95	0.99	0.11	-0.11	0.30	0.63
	LSTM MT	0.89	0.95	0.97	0.04	-0.19	0.07	0.56
	CNN-LSTM MT	0.82	0.94	1.09	0.22	-0.06	0.53	0.55
Lancang	LSTM	0.92	0.97	1.11	0.07	0.04	0.01	0.07
	CNN-LSTM	0.89	0.95	0.97	0.00	-0.11	0.02	0.18
	LSTM MT	0.93	0.96	0.99	-0.02	-0.04	-0.03	0.13
	CNN-LSTM MT	0.82	0.94	1.14	0.17	-0.01	0.09	0.32

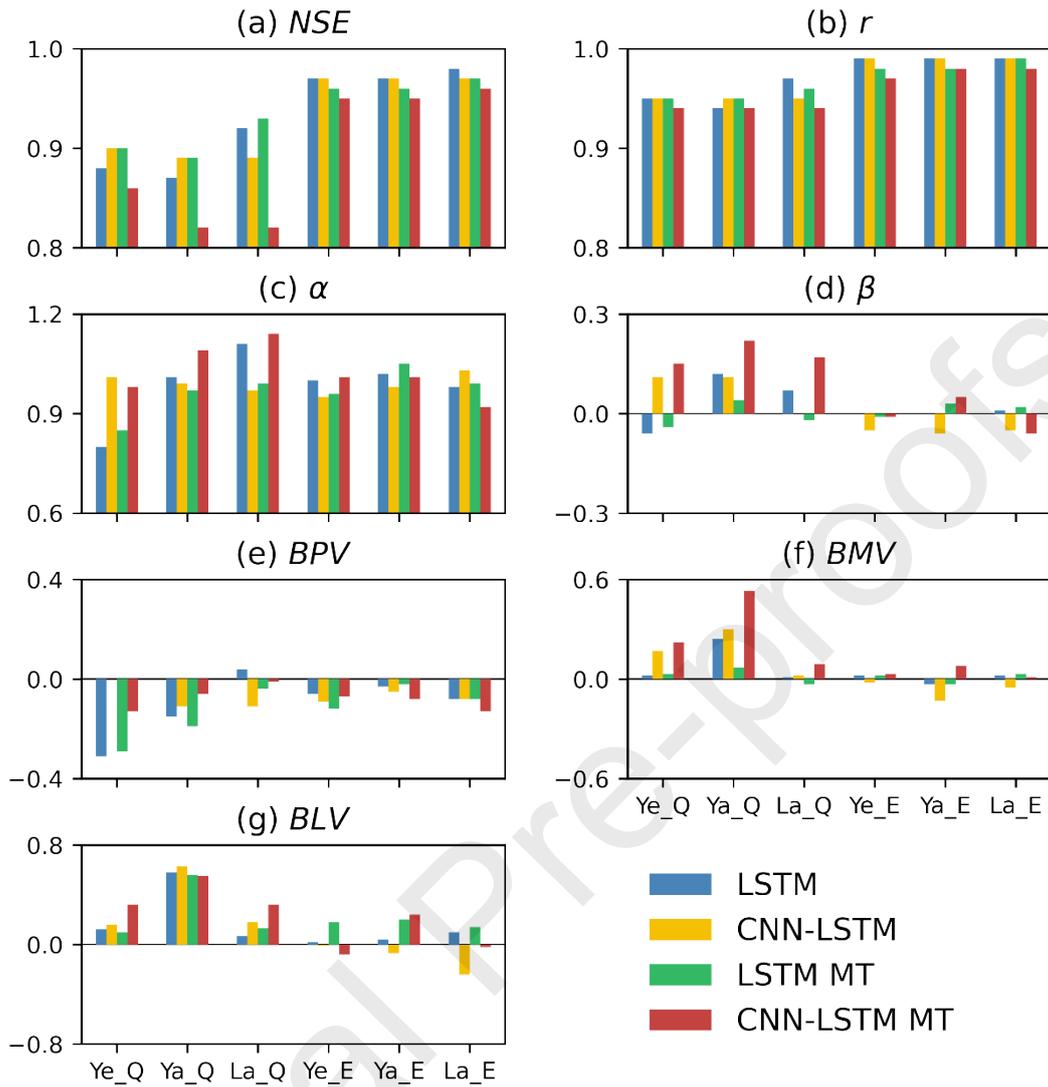


Figure 4. The evaluation metrics results of different DL-based hydrological models. “LSTM” (“CNN-LSTM”) and “LSTM MT” (“CNN-LSTM MT”) stand for LSTM (CNN-LSTM) ST and MT models, respectively. “Ye_Q”, “Ye_E”, “Ya_Q”, “Ya_E”, “La_Q”, and “La_E” denote the Q and E_a simulation results in the Yellow, Yangtze, and Lancang, respectively.

318 **Table 4.** The evaluation metrics of LSTM and CNN-LSTM hydrological models for
 319 E_a simulation. LSTM (MT) and CNN-LSTM (MT) represent the LSTM and CNN-
 320 LSTM E_a (multi-task) models, respectively.

Basins	Models	NSE	r	α	β	BPV	BMV	BLV
--------	--------	-------	-----	----------	---------	-------	-------	-------

Yellow	LSTM	0.97	0.99	1.00	0.00	-0.06	0.02	0.02
	CNN-LSTM	0.97	0.99	0.95	-0.05	-0.09	-0.02	-0.01
	LSTM MT	0.96	0.98	0.96	-0.01	-0.12	0.02	0.18
	CNN-LSTM MT	0.95	0.97	1.01	-0.01	-0.07	0.03	-0.08
Yangtze	LSTM	0.97	0.99	1.02	0.00	-0.03	-0.03	0.04
	CNN-LSTM	0.97	0.99	0.98	-0.06	-0.05	-0.13	-0.07
	LSTM MT	0.96	0.98	1.05	0.03	-0.02	-0.03	0.20
	CNN-LSTM MT	0.95	0.98	1.01	0.05	-0.08	0.08	0.24
Lancang	LSTM	0.98	0.99	0.98	0.01	-0.08	0.02	0.10
	CNN-LSTM	0.97	0.99	1.03	-0.05	-0.08	-0.05	-0.24
	LSTM MT	0.97	0.99	0.99	0.02	-0.08	0.03	0.14
	CNN-LSTM MT	0.96	0.98	0.92	-0.06	-0.13	0.01	-0.02

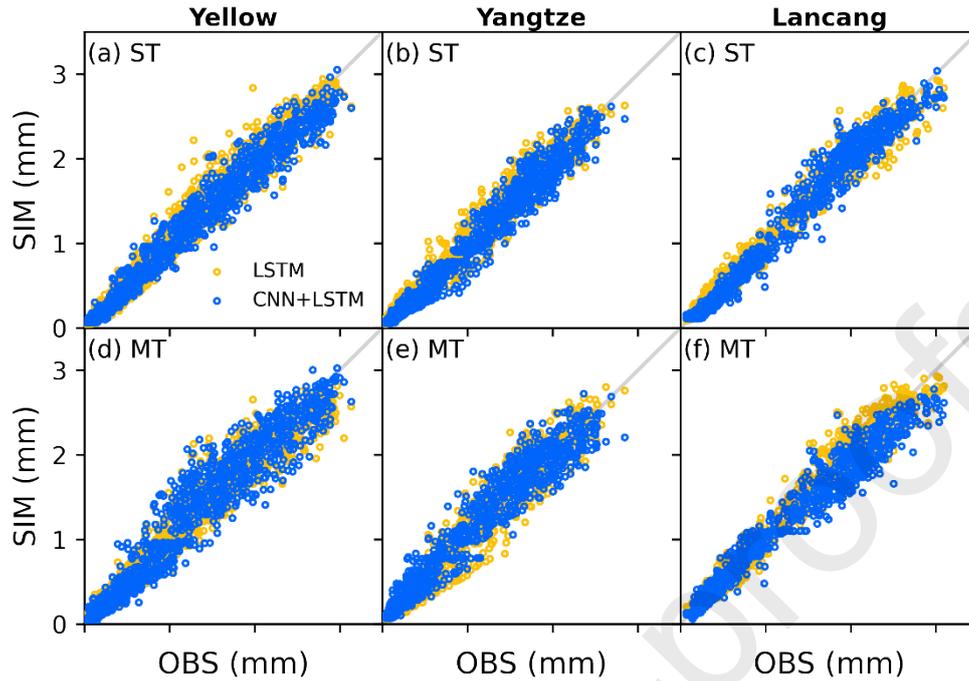


Figure 5. The comparison in daily E_a between simulations (SIM in y axis) and GLEAM (OBS in x axis) under two learning modes (ST and MT) in the test period at three study basins.

321 3.3.2 The capacity of 2-D CNN in extracting spatial characteristics

322 To evaluate the capacity of 2-D CNN in extracting spatial characteristics, we
 323 developed LSTM Q and E_a models as benchmarking baselines. High corrections
 324 ($NSEs$ more than 0.87 of Q and 0.97 of E_a ; Table 3-4) indicate the good
 325 performance of LSTM models in hydrological simulation. By comparing the
 326 performance of LSTM and CNN-LSTM Q models, it is found that CNN-LSTM
 327 models outperform LSTM models in overall Q simulation with higher $NSEs$ of Q
 328 in the Yellow and Yangtze (Table 3). Also, evaluation results of three biases metrics
 329 indicate that both LSTM and CNN-LSTM Q models underestimate the peak Q
 330 processes ($BPV < 0$) but overestimate the median and low Q processes
 331 ($BMV, BLV > 0$) in the Yellow and Yangtze. But the peak Q underestimation of
 332 CNN-LSTM models is mitigated compared with that of LSTM models, with BPV
 333 from -0.31 and -0.15 (LSTM) to 0.00 and -0.11 (CNN-LSTM) in the Yellow and
 334 Yangtze, respectively. In contrast, the overestimation of simulated median and low Q
 335 processes by CNN-LSTM Q models are stronger than that by LSTM Q models in
 336 all three study basins. We conclude that the introduction of 2-D CNN into LSTM Q
 337 models has a two-sided effect: it can enhance the overall model performance as well

338 as the capacity in simulating peaking Q processes, while may work slightly poorly in
339 modeling Q in the median and low ranges.

340 These effects can be explained by the different LSTM inputs. In the CNN-LSTM
341 models, the inputs to LSTM are the feature vectors (size is 1×64 at each time step;
342 Table A3) extracted from 2-D spatial meteorological data by 2-D CNN; whereas those
343 are the spatially-averaged P and T data within study basins in the LSTM-only
344 models. Based on the remarkable performance of both LSTM and CNN-LSTM
345 models, we assume that the LSTM inputs of CNN-LSTM models (i.e., CNN outputs)
346 contain the information about basin spatially-averaged P and T data (discussed later
347 in Sec. 4.2.1). Besides, we infer that more information hidden in the CNN outputs,
348 capturing the peak runoff features, contribute to enhancing the downstream LSTM
349 performance; while remaining redundant information might hurt downstream LSTM
350 performance. Noting that there is one exception: in the Lancang, CNN-LSTM Q
351 model does not outperform LSTM models in terms of overall and high Q : this could
352 be due to different hydrological characteristics of the testing period and shorter
353 training dataset (Table 2). All three years of the testing dataset in the Lancang can be
354 categorized as normal or dry years without extreme Q processes (e.g., 2012 in the
355 Yellow)—that might not benefit from the capacity of 2-D CNN in modeling peaking
356 Q processes. Besides, we also infer that shorter training dataset might not contain
357 enough information to train the optimal parameters for the CNN-LSTM model, whose
358 training parameters are larger than the LSTM-only model. It causes the CNN-LSTM
359 model does not outperform LSTM model in the Lancang. On the other hand, both
360 LSTM and CNN-LSTM E_a models perform well without significant differences
361 between two models in all three study basins (Table 4; Figure 4-5). It demonstrates
362 that LSTM is a remarkable approach for average E_a of the basin modeling and
363 introducing 2-D CNN has little impact on it, further suggesting that average P and T
364 play vital roles in E_a modeling and CNN outputs might contain the information
365 about it.

366 3.3.3 The effect of multi-task learning

367 To investigate the applicability of MT learning in training LSTM and CNN-LSTM
368 models, we configured them with Q and E_a as the simultaneous learning targets in
369 addition to those single-task (ST) experiments targeting at only Q or E_a . The
370 simulations by MT models agree well with observations in all three study basins with
371 favorable $NSEs$ of Q (0.89 by LSTM-MT and 0.82 by CNN-LSTM-MT) (Figure
372 6). Besides, both LSTM and CNN-LSTM MT models achieve remarkable
373 performance in predicting E_a in all three study basins (Figure 5; Table 4). These
374 results demonstrate that both LSTM and CNN-LSTM models are capable of
375 simulating multiple hydrological processes simultaneously.

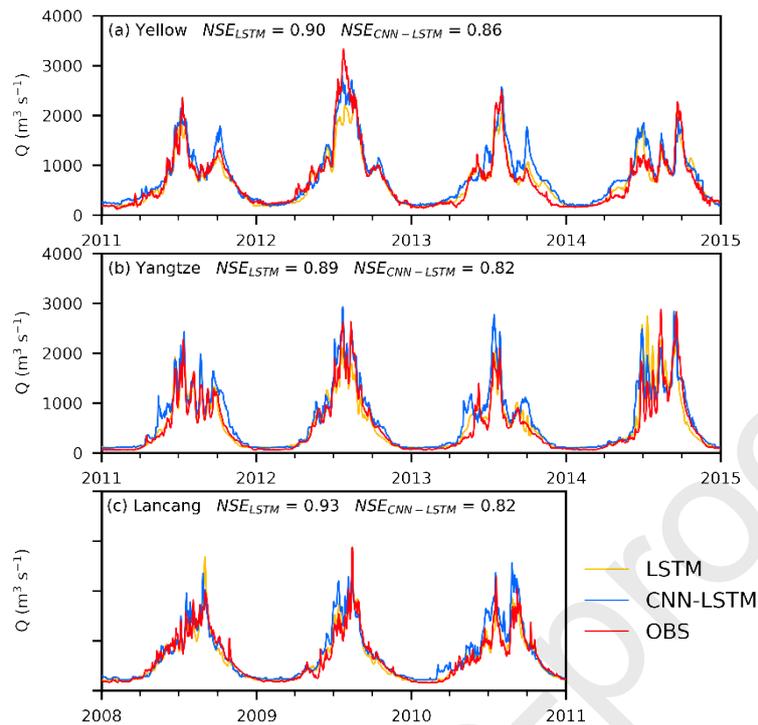


Figure 6. The comparison of simulated (MT models) and observed Q processes in the test period at three hydrological stations.

376 Furthermore, compared with LSTM Q models, the LSTM MT ones achieve higher
 377 $NSEs$ and lower total volume biases (β), suggesting that introducing E_a as
 378 additional training target can enhance the Q performance of LSTM models by
 379 reducing the total volume biases. It is consistent with the effect of introducing the E_a
 380 process in calibration for physical-based hydrological models (Herman et al. 2018,
 381 Nesru et al. 2020). It further motivates us to assume that the internal states of LSTM
 382 hydrological models, similar to physical-based hydrological models, might encompass
 383 multiple untrained physical hydrological processes. Thus, we put forward the
 384 hypothesis that E_a (Q) process information can be easily reconstructed from internal
 385 states of trained LSTM Q (E_a) models without the aforementioned complicated
 386 training procedures. This hypothesis will be tested in Sec. 4.2.2.

387 In contrast, the Q performance of CNN-LSTM MT models declines significantly
 388 compared with their ST counterparts in all three study basins. Significant performance
 389 differences between ST and MT might be explained by the variant LSTM inputs in
 390 CNN-LSTM models with different training references, compared with the invariant
 391 LSTM inputs in LSTM models (Sec. 2.3). Besides, the redundant information
 392 contained in the CNN outputs might make it difficult to simulate multiple

393 hydrological processes. On the other hand, differences in E_a performance between
 394 ST and MT mode of both LSTM and CNN-LSTM are not significant, indicating that
 395 MT learning in the LSTM and CNN-LSTM models has minimal effect on E_a
 396 modeling.

397 4. Physical interpretability of DL hydrological models

398 This section further explores the physical interpretability of CNN and LSTM models
 399 by using the LR model (Sec. 2.2) to look into potential hydrological and
 400 meteorological information CNN and LSTM may capture (Figure 7).

401 4.1 Experiment design

402 According to the above discussions, we design two experiments to test two
 403 hypotheses as follows:

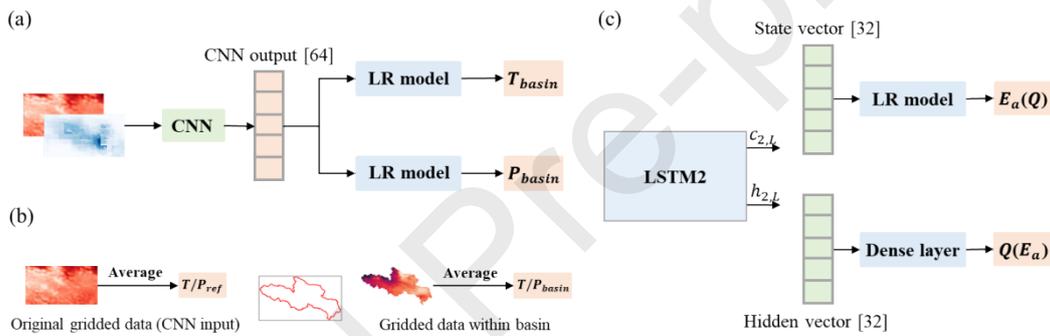


Figure 7. An overview of the physical interpretability exploration based on the LR model in this study. (a) LR models are used to map CNN output to spatially-averaged meteorological data within the basin; (b) The difference between original gridded data and gridded data within the basin; (c) LR models map the state vector of LSTM $Q(E_a)$ models to untrained $E_a(Q)$ processes.

404 (1) The CNN of trained CNN-LSTM models can capture spatially-averaged
 405 meteorological data within basin (T_{basin}/P_{basin}) from gridded meteorological
 406 data (CNN input; Figure 7a-b). To test this hypothesis, we construct LR
 407 models to map the CNN output of three CNN-LSTM (Q , E_a , and MT) models
 408 (Figure 7a) to spatially-averaged precipitation and temperature data within
 409 basins (P_{basin} and T_{basin}) and then use the correlation coefficient r between
 410 LR model output and spatially-averaged meteorological data within basins
 411 (namely r_{CNN}) as indicator to examine their correlation (Figure 7a). In
 412 addition, we calculate the r of spatial averages of meteorological data within

413 the basin and those of the original gridded data (r_{ref}) as the reference indicator
 414 (Figure 7b). By comparing r_{CNN} and r_{ref} , we can thus infer if the CNN can
 415 extract basin-scale characteristics of meteorological input data or not.

416 (2) E_a (Q) process information can be reconstructed from internal states of
 417 trained LSTM Q (E_a) models without the aforementioned complicated
 418 training procedures (Figure 7c). Similar as Lees et al. (2022), LR models are
 419 constructed to map the cell states vector—the memory of LSTM E_a (Q)
 420 models—to the untrained Q (E_a) processes and examine the resultant NSE to
 421 justify if LSTM is capable of inferring untrained hydrological information. For
 422 this, we obtain 6 sets of NSE values: 2 relations ($E_a \rightarrow Q$ and $Q \rightarrow E_a$) \times 3
 423 basins (Yellow, Yangtze, and Lancang).

424 4.2 The physical interpretability results

425 4.2.1 The physical interpretability of CNN outputs

426 As the meteorological data within a studied basin is essentially a subset of the original
 427 gridded dataset (Figure 7b), it is expected the spatial average of the former datasets
 428 (i.e., P_{basin} and T_{basin}) can be highly correlated to those of the latter (i.e., P_{ref} and
 429 T_{ref}), in particular when their spatial extents are comparable (e.g. the Yangtze river
 430 basin in Figure 2c). And our estimation does suggest high P and T correlation for
 431 all three basins ($r_{ref} > 0.97$) except the lower P correlation for the Lancang (r_{ref}
 432 < 0.80 ; Table 5). Thus, we select the Lancang basin as the study area to examine the
 433 capacity of CNN in inferring basin-specific information from the original gridded
 434 dataset. Our results indicate that the spatial averages of P and T within the Lancang
 435 can be extracted from the CNN outputs in high fidelity: r_{CNN} of all configurations
 436 are much larger than r_{ref} for P (cf. > 0.87 vs. -0.50 ; Table 5) and close to 1 for
 437 T , which supports our first hypothesis. Furthermore, it is also interesting to see the
 438 difference in r_{CNN} between different configurations: for P , a lower r_{CNN} is
 439 produced by the evapotranspiration-targeted model, suggesting those runoff-targeted
 440 models perform better in extracting precipitation-related information; while a reverse
 441 pattern is found for T . Such results are in line with the common hydrological
 442 understanding that P has a more significant impact on Q while T is a strong
 443 control of E_a .

444 **Table 5.** The r_{CNN} (three CNN-LSTM models: Q : CNN-Q, E_a : CNN-E, $Q + E_a$:
 445 CNN-Q+E) and r_{ref} of P and T in the Lancang.

r_{CNN}

r_{ref}

	CNN_Q	CNN_E	CNN_Q+E	
P	0.96	0.94	0.97	0.79
T	0.98	1.00	1.00	1.00

446 4.2.2 The physical interpretability of LSTM cell states

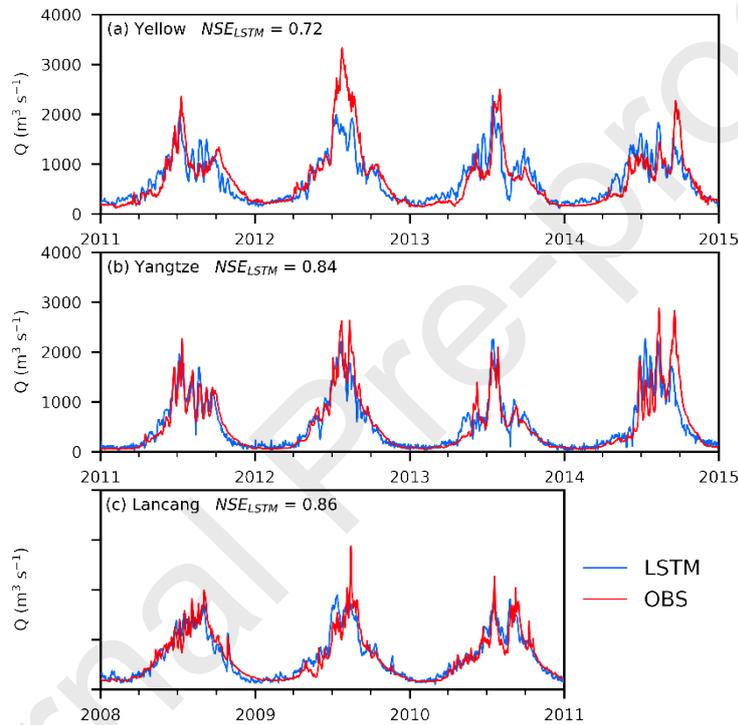


Figure 8. The comparison of simulated (extracted from the cell states of LSTM E_a models) and observed Q processes in the test period at three hydrological stations.

447 High correlation is found between the LSTM-inferred and observed Q ($NSEs$ of
 448 0.72, 0.84, and 0.86 in the Yellow, Yangtze, and Lancang; Figure 8), suggesting the
 449 remarkable capacity of LSTM in inferring the Q -related processes solely from E_a .
 450 Likewise, the inferred E_a processes also agree favorably with GLEAM E_a data with
 451 $NSEs$ more than 0.93 in all three study basins (Figure 9). Such results support our
 452 second hypothesis that LSTM Q (E_a) models contain information about untrained
 453 E_a (Q) processes without prior training procedures.

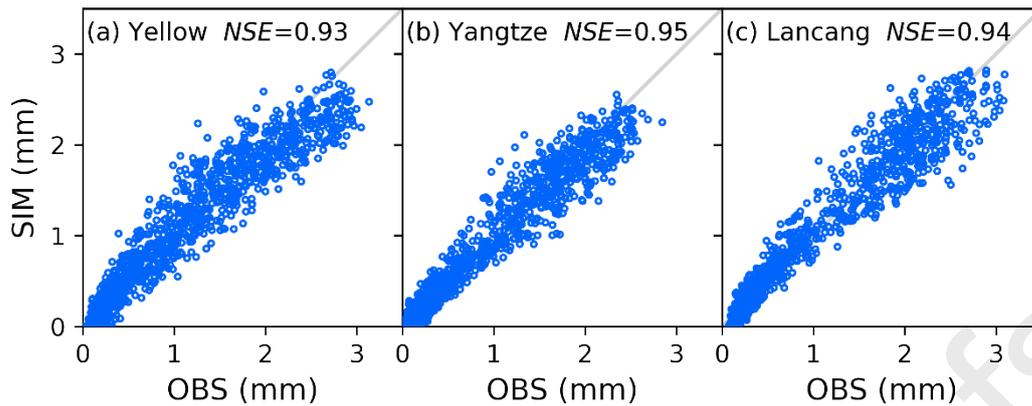


Figure 9. The comparison of simulated (extracted from the cell states of LSTM Q models) and observed E_a processes in the test period at three hydrological stations.

454 5. Conclusion and limitation

455 In this study, we develop an integrated spatiotemporal DL-based model by coupling
 456 CNN and LSTM for hydrological simulation and evaluate it in the source regions of
 457 the Yellow River, the Yangtze River, and the Lancang River. Besides, we employ a
 458 simple linear regression method to explore the physical interpretability of CNN and
 459 LSTM in DL-based hydrological models to improve our understanding of DLM-
 460 inferred hydrological processes. Our main findings are as follows:

- 461 (1) Both LSTM and CNN-LSTM models perform favorably with $NSEs$ more
 462 than 0.87 (Q) and 0.97 (E_a) in all three study basins. The CNN-LSTM Q
 463 models achieve better performance than LSTM ones in modeling Q .
 464 Therefore, we recommend the CNN-LSTM models for hydrological modeling
 465 and flooding forecasts in large mountainous basins.
- 466 (2) Both LSTM and CNN-LSTM multi-task models work well for Q (NSE
 467 $s > 0.82$) and E_a ($NSEs > 0.95$) simulation and LSTM multi-task models
 468 outperform LSTM Q models in all three study basins. This finding
 469 demonstrates that introducing E_a as an additional training target in the LSTM
 470 Q model can enhance the model performance and further research can be
 471 conducted to evaluate the potential of LSTM models for more and even whole
 472 hydrological modeling in the future.
- 473 (3) The internal cells of CNN and LSTM contain the hydrological and
 474 meteorological information consistent with our knowledge. CNN outputs in
 475 CNN-LSTM models can capture the basin-specific characteristics from 2-D
 476 gridded meteorological data of larger domains, while the internal cells of

477 LSTM Q (E_a) models may infer the E_a (Q) processes. These findings
478 advance the understanding of the physical interpretability of DL-based
479 hydrological models.

480 Being promising in predicting hydrological processes, this model has several
481 limitations. First, unlike the observed runoff data, the GLEAM dataset—as the E_a
482 training target—is of high accuracy but subject to some biases in the Tibetan Plateau
483 (Li et al. 2019a). It might lead trained E_a models to a slight deviation from the true
484 E_a process. Besides, this model is only evaluated in the three large mountainous
485 basins on the Tibetan Plateau due to the limitation of computational resources. Future
486 research will be carried out to evaluate the performance of proposed CNN-LSTM
487 model in more basins, develop CNN-LSTM regional hydrological models using more
488 underlying surface data, and explore more physical information hidden in the internal
489 cells of DL-based hydrological models.

490 **Appendix A: Architecture and parameters of LSTM and CNN-**
 491 **LSTM models**

492 This study developed 3 (basins: Yellow, Yangtze, Lancang) \times 2 (model: LSTM,
 493 CNN+LSTM) \times 3 (objective: runoff, evaporation, runoff + evaporation) = 18 DL-
 494 based hydrological models. The loss function is the mean-squared error (MSE), and
 495 the optimizer is Adaptive Moment Estimation (Adam). The batch size is 32, and the
 496 epoch is 400. The detailed parameters are as follows.

497 **Table A1.** The parameters of LSTM hydrological models

parameters	range of parameter values	final parameter values		
		Yello w	Yangtz e	Lancang
Hidden states	16, 32, 64, 128, 256	32	32	32
Length of the input sequence	10, 20, 30 \times n (n=1, 2, ..., 12)	90	180	120
Number of LSTM layer	1, 2	2	2	2

498

499 **Table A2.** The parameters of LSTM models in CNN-LSTM hydrological models

parameters	range of parameter values	final parameter values		
		Yello w	Yangtz e	Lancang
Hidden states	16, 32, 64, 128, 256	32	32	32
Length of the input sequence	10, 20, 30 \times n (n=1, 2, ..., 12)	60	120	90

Number of LSTM layer	1, 2	2	2	2
----------------------	------	---	---	---

500

501 **Table A3.** The parameters of CNN models in CNN-LSTM hydrological models

Yellow		Yangtze		Lancang	
size	Architectures	size	Architectures	size	architectures
41×77	<i>conv</i> , 5×5, s=2, 8	34×71	<i>conv</i> , 5×5, s=2, 8	84×59	<i>conv</i> , 5×5, s=2, 8
19×37	<i>conv</i> , 5×5, s=2, 16	16×35	<i>conv</i> , 5×5, s=2, 16	40×28	<i>conv</i> , 5×5, s=2, 16
9×17	<i>conv</i> , 5×5, s=2×3, 32	7×17	<i>conv</i> , 5×5, s=2×3, 32	19×12	<i>conv</i> , 5×5, s=3×2, 32
3×5	<i>conv</i> , 3×5, s=1, 64	2×5	<i>conv</i> , 2×5, s=1, 64	5×4	<i>conv</i> , 5×4, s=1, 64
1×1		1×1		1×1	

502 **Appendix B: The simulated daily and monthly runoff results of three**
503 **study basins in other relevant studies**

504 To compare the proposed CNN-LSTM model with traditional hydrological models,
505 we collected simulated daily and monthly runoff results using hydrological models in
506 relevant studies (Table B1). Although the simulation periods and model inputs are
507 different from this study, the higher *NSE* results can demonstrate that the CNN-
508 LSTM models outperform traditional hydrological models in runoff simulation.

509 **Table B1.** The simulated daily and monthly runoff *NSE* results of three study basins in
510 other relevant studies

Basins	Models	Simulation periods	<i>NSE</i> results	Reference
--------	--------	--------------------	--------------------	-----------

Yellow	SWAT	2014-2018	0.71	(Xie et al. 2020)
Yellow	SPHY	2000-2016	0.76	(Zhang et al. 2022)
Yangtze	WEB-DHM-SF	1987-2016	0.62	(Qi et al. 2019)
Yangtze	SWAT	2001-2016	0.75 (monthly)	(Ahmed et al. 2022)
Lancang	SWAT	1981-2015	0.71 (monthly)	(Li et al. 2021b)

511

512 *Code and data availability.* Precipitation, air temperature, and actual evaporation data
513 used in this research study are openly available and the sources are mentioned in Sec.
514 3.1.2. Model outputs and code are available by request to the corresponding author.
515 The observed runoff data are not publicly available for legal/ethical reasons.

516

517 **Acknowledgments**

518 This work was funded by National Natural Science Foundation of China (92047301)
519 and National Key Research and Development Project of China (2018YFA0606002).

520

521 **References**

- 522 Ahmed, N., Wang, G., Booij, M.J., et al. (2022). Separation of the impact of
523 landuse/landcover change and climate change on runoff in the upstream area
524 of the Yangtze River, China. *Water resources management* 36(1), 181-201.
- 525 Ali, A.S.A., Ebrahimi, S., Ashiq, M.M., et al. (2022). CNN-Bi LSTM neural network
526 for simulating groundwater level. *Environ Eng* 8, 1-7.
- 527 Arrieta, A.B., Díaz-Rodríguez, N., Del Ser, J., et al. (2020). Explainable Artificial
528 Intelligence (XAI): Concepts, taxonomies, opportunities and challenges
529 toward responsible AI. *Information Fusion* 58, 82-115.
- 530 Barzegar, R., Aalami, M.T. and Adamowski, J. (2020). Short-term water quality
531 variable prediction using a hybrid CNN–LSTM deep learning model.
532 *Stochastic Environmental Research and Risk Assessment* 34(2), 415-433.
- 533 Beck, H.E., van Dijk, A.I.J.M., Levizzani, V., et al. (2017). MSWEP: 3-hourly 0.25°
534 global gridded precipitation (1979–2015) by merging gauge, satellite, and
535 reanalysis data. *Hydrology and Earth System Sciences* 21(1), 589-615.
- 536 Beck, H.E., Wood, E.F., Pan, M., et al. (2019). MSWEP V2 Global 3-Hourly 0.1°
537 Precipitation: Methodology and Quantitative Assessment. *Bulletin of the*
538 *American Meteorological Society* 100(3), 473-500.
- 539 Bengio, Y., Simard, P. and Frasconi, P. (1994). Learning long-term dependencies
540 with gradient descent is difficult. *IEEE transactions on neural networks* 5(2),
541 157-166.
- 542 Blöschl, G., Bierkens, M.F.P., Chambel, A., et al. (2019). Twenty-three unsolved
543 problems in hydrology (UPH) – a community perspective. *Hydrological*
544 *Sciences Journal* 64(10), 1141-1158.
- 545 Blöschl, G. and Sivapalan, M. (1995). Scale issues in hydrological modelling: a
546 review. *Hydrological Processes* 9(3-4), 251-290.
- 547 Caruana, R. (1997). Multitask learning. *Machine learning* 28(1), 41-75.
- 548 Demirel, M.C., Venancio, A. and Kahya, E. (2009). Flow forecast by SWAT model
549 and ANN in Pracana basin, Portugal. *Advances in Engineering Software*
550 40(7), 467-473.
- 551 Duan, S. and Ullrich, P. (2021). A comprehensive investigation of machine learning
552 models for estimating daily snow water equivalent over the Western US. *Earth*

- 553 and Space Science Open Archive.
- 554 Feng, D., Fang, K. and Shen, C. (2020). Enhancing Streamflow Forecast and
555 Extracting Insights Using Long-Short Term Memory Networks With Data
556 Integration at Continental Scales. *Water Resources Research* 56(9),
557 e2019WR026793.
- 558 Frame, J.M., Kratzert, F., Raney, A., et al. (2021). Post-Processing the National Water
559 Model with Long Short-Term Memory Networks for Streamflow Predictions
560 and Model Diagnostics. *Journal of the American Water Resources Association*
561 57(6), 885-905.
- 562 Gupta, H.V., Kling, H., Yilmaz, K.K., et al. (2009). Decomposition of the mean
563 squared error and NSE performance criteria: Implications for improving
564 hydrological modelling. *Journal of Hydrology* 377(1-2), 80-91.
- 565 Gupta, H.V., Wagener, T. and Liu, Y. (2008). Reconciling theory with observations:
566 elements of a diagnostic approach to model evaluation. *Hydrological*
567 *Processes: An International Journal* 22(18), 3802-3813.
- 568 Herman, M.R., Nejadhashemi, A.P., Abouali, M., et al. (2018). Evaluating the role of
569 evapotranspiration remote sensing data in improving hydrological modeling
570 predictability. *Journal of Hydrology* 556, 39-49.
- 571 Hersbach, H., Bell, B., Berrisford, P., et al. (2020). The ERA5 global reanalysis.
572 *Quarterly Journal of the Royal Meteorological Society* 146(730), 1999-2049.
- 573 Hewitt, J. and Liang, P. (2019). Designing and Interpreting Probes with Control
574 Tasks. *Proceedings of the 2019 Con.*
- 575 Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural*
576 *computation* 9(8), 1735-1780.
- 577 Hrachowitz, M., Savenije, H.H.G., Blöschl, G., et al. (2013). A decade of Predictions
578 in Ungauged Basins (PUB)—a review. *Hydrological Sciences Journal* 58(6),
579 1198-1255.
- 580 Hsu, K.I., Gupta, H.V. and Sorooshian, S.J.W.r.r. (1995). Artificial neural network
581 modeling of the rainfall-runoff process. *Water Resources Research* 31(10),
582 2517-2530.
- 583 Huss, M., Bookhagen, B., Huggel, C., et al. (2017). Toward mountains without
584 permanent snow and ice. *Earths Future* 5(5), 418-435.

- 585 Immerzeel, W.W., Van Beek, L.P. and Bierkens, M.F. (2010). Climate change will
586 affect the Asian water towers. *Science* 328(5984), 1382-1385.
- 587 J.E.Nash and J.V.Sutcliffe (1970). River flow forecasting through conceptual models
588 part I — A discussion of principles. *Journal of Hydrology* 10(3), 282-290.
- 589 Jiang, S., Zheng, Y. and Solomatine, D. (2020). Improving AI System Awareness of
590 Geoscience Knowledge: Symbiotic Integration of Physical Approaches and
591 Deep Learning. *Geophysical Research Letters* 47(13).
- 592 Jiang, S., Zheng, Y., Wang, C., et al. (2022). Uncovering Flooding Mechanisms
593 Across the Contiguous United States Through Interpretive Deep Learning on
594 Representative Catchments. *Water Resources Research* 58(1),
595 e2021WR030185.
- 596 Khan, M.Y.A., Tian, F., Hasan, F., et al. (2019). Artificial neural network simulation
597 for prediction of suspended sediment concentration in the River Ramganga,
598 Ganges Basin, India. *International Journal of Sediment Research* 34(2), 95-
599 107.
- 600 Khatakho, R., Talchabhadel, R. and Thapa, B.R. (2021). Evaluation of different
601 precipitation inputs on streamflow simulation in Himalayan River basin.
602 *Journal of Hydrology* 599.
- 603 Kratzert, F., Klotz, D., Brenner, C., et al. (2018). Rainfall–runoff modelling using
604 Long Short-Term Memory (LSTM) networks. *Hydrology and Earth System
605 Sciences* 22(11), 6005-6022.
- 606 Kratzert, F., Klotz, D., Shalev, G., et al. (2019). Towards learning universal, regional,
607 and local hydrological behaviors via machine learning applied to large-sample
608 datasets. *Hydrology and Earth System Sciences* 23(12), 5089-5110.
- 609 LeCun, Y., Bottou, L., Bengio, Y., et al. (1998). Gradient-based learning applied to
610 document recognition. *Proceedings of the IEEE* 86(11), 2278-2324.
- 611 Lees, T., Buechel, M., Anderson, B., et al. (2021). Benchmarking data-driven
612 rainfall–runoff models in Great Britain: a comparison of long short-term
613 memory (LSTM)-based models with four lumped conceptual models.
614 *Hydrology and Earth System Sciences* 25(10), 5517-5534.
- 615 Lees, T., Reece, S., Kratzert, F., et al. (2022). Hydrological concept formation inside
616 long short-term memory (LSTM) networks. *Hydrology Earth System Sciences*
617 26(12), 3079-3101.

- 618 Li, B., Zhou, X., Ni, G., et al. (2021a). A multi-factor integrated method of
619 calculation unit delineation for hydrological modeling in large mountainous
620 basins. *Journal of Hydrology* 597, 126180.
- 621 Li, R., Huang, H., Yu, G., et al. (2021b). Contributions of climatic variation and
622 human activities to streamflow changes in the Lancang-Mekong River Basin.
623 *Resources Science* 43(12), 2428-2441.
- 624 Li, R., Sun, T., Tian, F., et al. (2022). SHAFTS (v2022. 3): a deep-learning-based
625 Python package for Simultaneous extraction of building Height And
626 Footprint from Sentinel Imagery. *Geoscientific Model Development*
627 *Discussions*, 1-42.
- 628 Li, R., Sun, T., Tian, F., et al. (2023). SHAFTS (v2022.3): a deep-learning-based
629 Python package for simultaneous extraction of building height and footprint
630 from sentinel imagery. *Geoscientific Model Development* 16(2), 751-778.
- 631 Li, X., Long, D., Han, Z., et al. (2019a). Evapotranspiration estimation for Tibetan
632 plateau headwaters using conjoint terrestrial and atmospheric water balances
633 and multisource remote sensing. *Water Resources Research* 55(11), 8608-
634 8630.
- 635 Li, Z., Feng, Q., Li, Z., et al. (2019b). Climate background, fact and hydrological
636 effect of multiphase water transformation in cold regions of the Western
637 China: A review. *Earth-Science Reviews* 190, 33-57.
- 638 Liu, Y., Zhang, T., Kang, A., et al. (2021). Research on Runoff Simulations Using
639 Deep-Learning Methods. *Sustainability* 13(3), 1336.
- 640 Martens, B., Miralles, D.G., Lievens, H., et al. (2017). GLEAM v3: satellite-based
641 land evaporation and root-zone soil moisture. *Geoscientific Model*
642 *Development* 10(5), 1903-1925.
- 643 Miao, Q., Pan, B., Wang, H., et al. (2019). Improving Monsoon Precipitation
644 Prediction Using Combined Convolutional and Long Short Term Memory
645 Neural Network. *Water* 11(5), 977.
- 646 Miralles, D.G., Holmes, T.R.H., De Jeu, R.A.M., et al. (2011). Global land-surface
647 evaporation estimated from satellite-based observations. *Hydrology and Earth*
648 *System Sciences* 15(2), 453-469.
- 649 Moriasi, D.N., Arnold, J.G., Van Liew, M.W., et al. (2007). Model evaluation
650 guidelines for systematic quantification of accuracy in watershed simulations.
651 *Transactions of the Asabe* 50(3), 885-900.

- 652 Nan, Y., He, Z., Tian, F., et al. (2021). Can we use precipitation isotope outputs of
653 isotopic general circulation models to improve hydrological modeling in large
654 mountainous catchments on the Tibetan Plateau? *Hydrology and Earth System*
655 *Sciences* 25(12), 6151-6172.
- 656 Nearing, G.S., Kratzert, F., Sampson, A.K., et al. (2021). What Role Does
657 Hydrological Science Play in the Age of Machine Learning? *Water Resources*
658 *Research* 57(3), e2020WR028091.
- 659 Nesru, M., Shetty, A. and Nagaraj, M.K. (2020). Multi-variable calibration of
660 hydrological model in the upper Omo-Gibe basin, Ethiopia. *Acta Geophysica*
661 68(2), 537-551.
- 662 Nourani, V., Khodkar, K. and Gebremichael, M. (2022). Uncertainty assessment of
663 LSTM based groundwater level predictions. *Hydrological Sciences Journal*
664 67(5), 773-790.
- 665 Qi, J., Wang, L., Zhou, J., et al. (2019). Coupled snow and frozen ground physics
666 improves cold region hydrological simulations: an evaluation at the upper
667 Yangtze River Basin (Tibetan Plateau). *Journal of Geophysical Research:*
668 *Atmospheres* 124(23), 12985-13004.
- 669 Reichstein, M., Camps-Valls, G., Stevens, B., et al. (2019). Deep learning and process
670 understanding for data-driven Earth system science. *Nature* 566(7743), 195-
671 204.
- 672 Sadeghi Tabas, S. and Samadi, S. (2022). Variational Bayesian dropout with a
673 Gaussian prior for recurrent neural networks application in rainfall–runoff
674 modeling. *Environmental Research Letters* 17(6).
- 675 Schaner, N., Voisin, N., Nijssen, B., et al. (2012). The contribution of glacier melt to
676 streamflow. *Environmental Research Letters* 7(3), 034029.
- 677 Sherstinsky, A. (2020). Fundamentals of Recurrent Neural Network (RNN) and Long
678 Short-Term Memory (LSTM) network. *Physica D: Nonlinear Phenomena* 404.
- 679 Shi, X., Chen, Z., Wang, H., et al. (2015). Convolutional LSTM network: A machine
680 learning approach for precipitation nowcasting. *Advances in neural*
681 *information processing systems* 28.
- 682 Sood, A. and Smakhtin, V. (2015). Global hydrological models: a review.
683 *Hydrological Sciences Journal* 60(4), 549-565.
- 684 Sun, A.Y., Jiang, P., Mudunuru, M.K., et al. (2021). Explore Spatio-Temporal

- 685 Learning of Large Sample Hydrology Using Graph Neural Networks. *Water*
686 *Resources Research* 57(12), e2021WR030394.
- 687 Tarek, M., Brissette, F.P. and Arsenault, R. (2020). Evaluation of the ERA5
688 reanalysis as a potential reference dataset for hydrological modelling over
689 North America. *Hydrology and Earth System Sciences* 24(5), 2527-2544.
- 690 Wang, T., Zhao, Y., Xu, C., et al. (2021). Atmospheric dynamic constraints on
691 Tibetan Plateau freshwater under Paris climate targets. *Nature Climate Change*
692 11(3), 219-225.
- 693 Xie, P., Zhuo, L., Yang, X., et al. (2020). Spatial-temporal variations in blue and
694 green water resources, water footprints and water scarcities in a large river
695 basin: A case for the Yellow River basin. *Journal of Hydrology* 590, 125222.
- 696 Yang, S., Yang, D., Chen, J., et al. (2020). A physical process and machine learning
697 combined hydrological model for daily streamflow simulations of large
698 watersheds with limited observation data. *Journal of Hydrology* 590, 125206.
- 699 Yang, Y., Xiong, Q., Wu, C., et al. (2021). A study on water quality prediction by a
700 hybrid CNN-LSTM model with attention mechanism. *Environ Sci Pollut Res*
701 *Int* 28(39), 55129-55139.
- 702 Yilmaz, K.K., Gupta, H.V. and Wagener, T. (2008). A process-based diagnostic
703 approach to model evaluation: Application to the NWS distributed hydrologic
704 model. *Water Resources Research* 44(9), W09417.
- 705 Zhang, L., Su, F., Yang, D., et al. (2013). Discharge regime and simulation for the
706 upstream of major rivers over Tibetan Plateau. *Journal of Geophysical*
707 *Research: Atmospheres* 118(15), 8500-8518.
- 708 Zhang, T., Li, D. and Lu, X. (2022). Response of runoff components to climate
709 change in the source-region of the Yellow River on the Tibetan plateau.
710 *Hydrological Processes* 36(6), e14633.

711

712 **Table 1.** Basic facts of the three study basins. “DEM” represents Digital Elevation
713 Model.

Basins	Area	DEM range	Average annual P	Average annual Q
	(km ²)	(m)	(mm)	(10 ⁹ m ³)

Yellow	123,000	2,656-6,253	510	20
Yangtze	139,000	3,516-6,575	460	16
Lancang	91,000	1,243-6,334	830	32

714

Journal Pre-proofs

715 **Table 2.** The length of the training, evaluation, and testing periods in three study
716 basins.

Basins	Training period	Evaluation period	Testing period
Yellow	1983-2004	2006-2009	2011-2014
Yangtze	1983-2004	2006-2009	2011-2014
Lancang	1988-2004	2005-2007	2008-2010

717

718 **Table 3.** The evaluation metrics of LSTM and CNN-LSTM hydrological models for
 719 Q simulation. LSTM (MT) and CNN-LSTM (MT) represent the LSTM and CNN-
 720 LSTM Q (multi-task) models, respectively.

		NSE	r	α	β	BPV	BMV	BLV
Yellow	LSTM	0.88	0.95	0.80	-0.06	-0.31	0.02	0.12
	CNN-LSTM	0.90	0.95	1.01	0.11	-0.00	0.17	0.16
	LSTM MT	0.90	0.95	0.85	-0.04	-0.29	0.03	0.10
	CNN-LSTM MT	0.86	0.94	0.98	0.15	-0.13	0.22	0.32
Yangtze	LSTM	0.87	0.94	1.01	0.12	-0.15	0.24	0.58
	CNN-LSTM	0.89	0.95	0.99	0.11	-0.11	0.30	0.63
	LSTM MT	0.89	0.95	0.97	0.04	-0.19	0.07	0.56
	CNN-LSTM MT	0.82	0.94	1.09	0.22	-0.06	0.53	0.55
Lancang	LSTM	0.92	0.97	1.11	0.07	0.04	0.01	0.07
	CNN-LSTM	0.89	0.95	0.97	0.00	-0.11	0.02	0.18
	LSTM MT	0.93	0.96	0.99	-0.02	-0.04	-0.03	0.13
	CNN-LSTM MT	0.82	0.94	1.14	0.17	-0.01	0.09	0.32

721

722 **Table 4.** The evaluation metrics of LSTM and CNN-LSTM hydrological models for
 723 E_a simulation. LSTM (MT) and CNN-LSTM (MT) represent the LSTM and CNN-
 724 LSTM E_a (multi-task) models, respectively.

Basins	Models	NSE	r	α	β	BPV	BMV	BLV
Yellow	LSTM	0.97	0.99	1.00	0.00	-0.06	0.02	0.02
	CNN-LSTM	0.97	0.99	0.95	-0.05	-0.09	-0.02	-0.01
	LSTM MT	0.96	0.98	0.96	-0.01	-0.12	0.02	0.18
	CNN-LSTM MT	0.95	0.97	1.01	-0.01	-0.07	0.03	-0.08
Yangtze	LSTM	0.97	0.99	1.02	0.00	-0.03	-0.03	0.04
	CNN-LSTM	0.97	0.99	0.98	-0.06	-0.05	-0.13	-0.07
	LSTM MT	0.96	0.98	1.05	0.03	-0.02	-0.03	0.20
	CNN-LSTM MT	0.95	0.98	1.01	0.05	-0.08	0.08	0.24
Lancang	LSTM	0.98	0.99	0.98	0.01	-0.08	0.02	0.10
	CNN-LSTM	0.97	0.99	1.03	-0.05	-0.08	-0.05	-0.24
	LSTM MT	0.97	0.99	0.99	0.02	-0.08	0.03	0.14
	CNN-LSTM MT	0.96	0.98	0.92	-0.06	-0.13	0.01	-0.02

725

726 **Table 5.** The r_{CNN} (three CNN-LSTM models: Q : CNN-Q, E_a : CNN-E, $Q + E_a$:
727 CNN-Q+E) and r_{ref} of P and T in the Lancang.

	r_{CNN}			r_{ref}
	CNN_Q	CNN_E	CNN_Q+E	
P	0.96	0.94	0.97	0.79
T	0.98	1.00	1.00	1.00

728

729

730 **Declaration of interests**

731

732 The authors declare that they have no known competing financial interests or personal
733 relationships that could have appeared to influence the work reported in this paper.

734

735 The authors declare the following financial interests/personal relationships which may be
736 considered as potential competing interests:

737

738

739

740

741

742

743

744 Long short-term memory (LSTM) networks have demonstrated their excellent
745 capability in processing long-length temporal dynamics and have proven to be
746 effective in precipitation-runoff modeling. However, the current LSTM hydrological
747 models lack the incorporation of multi-task learning and spatial information, which
748 limits their ability to make full use of meteorological and hydrological data. To
749 address this issue, this study proposes a spatiotemporal deep-learning (DL)-based
750 hydrological model that couples the 2-Dimension convolutional neural network
751 (CNN) and LSTM and introduces actual evaporation (E_a) as an additional training

752 target. The proposed CNN-LSTM model is tested on three large mountainous basins
753 on the Tibetan Plateau, and the results are compared to those obtained from the
754 LSTM-only model. Additionally, a probe method is used to decipher the internal
755 embedding layers of the proposed DL models. The results indicate that both LSTM
756 and CNN-LSTM hydrological models perform well in simulating runoff (Q) and E_a ,
757 with Nash-Sutcliffe efficiency coefficients ($NSEs$) higher than 0.82 and 0.95,
758 respectively. The higher $NSEs$ suggest that introducing spatial information into
759 LSTM-only models can improve the overall and peak model performance. Moreover,
760 multi-task simulation with LSTM-only models shows better accuracy in the
761 estimation of Q volume and performance, with $NSEs$ increasing by approximately
762 0.02. The probe method also reveals that CNN can capture the basin-averaged
763 meteorological values in CNN-LSTM models, while LSTM Q (E_a) models contain
764 the information about the known E_a (Q) process. Overall, this study demonstrates the
765 value of spatial information and multi-task learning in LSTM hydrological modeling
766 and provides a perspective for interpreting the internal embedding layers of DL
767 models.

768

769