



1

2

3 **Semi-supervised learning approach to improve the predictability of**
4 **data-driven rainfall-runoff model in hydrological data-sparse**
5 **regions**

6

7

Sunghyun Yoon¹ and Kuk-Hyun Ahn²

8

9

10

June 2023

11

12

13

14

15

16

17

18

19

20 ¹Assistant Professor, Department of Artificial Intelligence, Kongju National University,
21 Cheon-an, South Korea; e-mail: syoon@kongju.ac.kr

22 ²Associate Professor, Department of Civil and Environmental Engineering, Kongju National
23 University, Cheon-an, South Korea; *Corresponding author*; e-mail: [ahnkukhyun@](mailto:ahnkukhyun@kongju.ac.kr)
24 [kongju.ac.kr](mailto:ahnkukhyun@kongju.ac.kr)



25 ABSTRACT

26

27 Numerous data-driven models have been introduced to establish reliable predictions in the
28 rainfall-runoff relationship. The majority of these models are trained using a supervised
29 learning (SL) approach, with paired observed samples of climate and streamflow data.
30 However, in practice, the availability of such paired observations is often constrained due to
31 sparse data from streamflow gauges worldwide, which typically covers only a few years. This
32 limited number of paired samples can significantly impede the learning ability of the data-
33 driven model. The semi-supervised learning approach, which is an emerging machine learning
34 paradigm that additionally incorporates unpaired samples, has the potential to be a highly
35 effective method for modeling rainfall-runoff relationships. In this study, we present a novel
36 semi-supervised learning-based framework for rainfall-runoff modeling. Our framework
37 introduces a unique loss function designed to handle two distinct types of samples, namely
38 paired and unpaired samples, effectively during the training process. To validate the
39 effectiveness of the proposed framework, we conducted an extensive set of experiments
40 employing a diverse range of designs, all of which utilized the LSTM network. The
41 experiments are based on 531 basins from the freely available CAMELS dataset, which spans
42 the entire continuous United States. Results indicate that the proposed framework show
43 significantly enhanced performance compared to the baseline models. Results also show that
44 the framework can serve as a viable alternative to the previously developed fully supervised
45 approaches. Lastly, we address potential avenues for enhancing the model and provide an
46 outline of our future research plans in this domain.

47 **Keywords:** Long short-term memory (LSTM); Semi-supervised learning; Data-sparse region;
48 Rainfall-runoff modeling; Unpaired samples;



49 **1. Introduction**

50 Rainfall-runoff modeling is an essential tool for urban planning, land use, flood and water
51 resource management (Nourani et al., 2009). It represents the hydrologic processes involved in
52 converting rainfall into runoff, making it one of the principal interests in hydrological sciences
53 (Beven, 2011; Sitterson et al., 2018). Over time, the modeling of the rainfall-runoff process has
54 evolved from physical-based models such as SHETRAN (Birkinshaw et al., 2010) and
55 VELMA (Mckane et al., 2014) to conceptual models such as Variable Infiltration Capacity
56 (VIC; Liang et al., 1994), Hydrologiska Byråns Vattenbalansavdelning (HBV; Seibert and Vis,
57 2012), and Sacramento Soil Moisture Accounting (SCA-SMA; Burnash et al., 1973). Data-
58 driven models have also been employed to depict the rainfall-runoff process, with recent
59 studies reporting their ability to outperform traditional models (Hoedt et al., 2021; Lees et al.,
60 2021; Reichstein et al., 2019; Xiang et al., 2020). In this paper, we aim to further improve the
61 predictive ability of the data-driven model, which is currently regarded as the state-of-the-art
62 in hydrologic prediction (Nearing et al., 2021; Shen et al., 2021).

63

64 Data-driven models leverage empirical relationships between target and independent variables,
65 offering the advantages of requiring low input, minimal effort for development and application,
66 and moderate computational resources (Abbott, 1999; Chen et al., 2018). Prominent data-
67 driven techniques include genetic programming (Chadalawada et al., 2020), support vector
68 machine (SVM) (Alquraish and Khadr, 2021), random forests (Booker and Woods, 2014), and
69 fuzzy logic (Bartoletti et al., 2018; Kothari and Gharde, 2015). Deep learning (DL) techniques
70 have also gained significant traction for their effectiveness (Roy et al., 2021; Taormina and
71 Chau, 2015; Van et al., 2020; Xie et al., 2021). One standout architecture, the long short-term



72 memory (LSTM; Hochreiter and Schmidhuber, 1997) network, has been specifically designed
73 to simulate time series by incorporating an inductive bias that preserves crucial temporal
74 information over extended periods (Hoedt et al., 2021).

75

76 LSTMs have been shown to provide a significant advantage over conventional hydrologic
77 models in rainfall-runoff modeling by a considerable margin (Kratzert et al., 2018; Lees et al.,
78 2021), even at hourly time scales (Gauch et al., 2021a), and for watersheds unseen by the LSTM
79 (Arsenault et al., 2023; Kratzert et al., 2019a). For instance, Kratzert et al. (2019b)
80 demonstrated that when using an LSTM to predict streamflow in 531 basins across the United
81 States (US), it outperformed several different hydrological benchmark models, including SAC-
82 SMA, VIC, and HBV models. In recent years, LSTMs have been utilized to (i) quantify the
83 predictive uncertainty (Klotz et al., 2022; Li et al., 2021), (ii) evaluate the suitability of
84 hydrologic projections under climate change (Wi and Steinschneider, 2022), and (iii) improve
85 the reliability of simulations in hydrologic models as post-processors (Frame et al., 2021; Hunt
86 et al., 2022). Notably, several studies have demonstrated exceptional LSTM performance,
87 especially in situations where abundant data are available (Anderson and Radic, 2021; Gauch
88 et al., 2021b; Lees et al., 2021).

89

90 However, acquiring hydrological records with comprehensive long-term coverage is often
91 unattainable in reality. Many regions worldwide face the challenge of limited streamflow gauge
92 networks, resulting in sparse data that typically spans only a few years (Bitew and
93 Gebremichael, 2011; Do et al., 2017). For example, Lee and Ahn (2022) have utilized a limited
94 number of only 27 streamflow gauges to investigate a national-scale hydrologic variability



95 across South Korea. Similarly, in the Tana River basin in Kenya, which covers an expansive
96 area of 95,000 km² and serves as a habitat for diverse wildlife species (TNC, 2015),
97 hydrological data from only 26 streamflow gauges spanning a five-year period (February 2015
98 to January 2020) are available (Leisher et al., 2016). The Global Runoff Data Center (GRDC)
99 provides daily streamflow observations for a significant portion of basins globally, with records
100 typically spanning less than three years in length. The availability of observed streamflow
101 records may also be problematic for data-rich regions due to several situations like gauges
102 installed in recent decades, discontinuation from budgetary constraints, or measuring
103 malfunctioning over an extended period of time (Ahn, 2021). Consequently, many data-driven
104 models have been applied in a local modeling context, wherein a model is trained using data
105 from one or a few basins (e.g., Bowes et al., 2019; Han et al., 2021; Ley et al., 2023; Liang et
106 al., 2018; Xu et al., 2022). The potential limitations associated with sparse records have been
107 discussed, and the need for corrective measures has been addressed (Beven, 2020; Shen, 2018).

108

109 In areas where streamflow records are scarce, longer historical climate data records often
110 remain available. However, current methods for training rainfall-runoff models in data-sparse
111 regions typically rely solely on paired recorded samples between climate and streamflow data,
112 known as labeled data in machine learning (ML) terminology. Nonetheless, valuable insights
113 can be gained by incorporating the remaining climate data, referred to as unlabeled data, to
114 improve model performance. In the field of ML, semi-supervised and unsupervised learning
115 are emerging paradigms that utilize unlabeled data to enhance model performance. While semi-
116 supervised learning combines both labeled and unlabeled data to improve performance,
117 unsupervised learning first pre-tunes with unlabeled data before fine-tuning with labeled data



118 (Chen et al., 2022; He et al., 2020). In particular, semi-supervised learning is increasingly
119 recognized as an effective approach with enhanced learning ability (Du et al., 2020; Levatić et
120 al., 2017; Zhou and Zhou, 2021).

121

122 While semi-supervised learning has gained popularity in various fields, its formal application
123 in the field of hydrology is still limited. In this study, we introduce a novel framework based
124 on semi-supervised learning for rainfall-runoff modeling, aiming to explore its potential
125 usefulness in structuring hydrological time series modeling problems. Specifically, we present
126 how LSTM models can enhance their predictive performance in regions with limited data,
127 thereby addressing the limitations associated with streamflow observations. In literature, we
128 found two previous studies focusing on improving the modeling performance in data-sparse
129 regions (Ma et al., 2021; Oruche et al., 2021). The approach we propose is notably distinct
130 from those studies in that we do not use any source datasets from other regions. Both of these
131 studies utilize transfer learning, a technique in which a pre-trained model from extensive
132 labeled data from other continents is used to transfer initial weights to a model. In this study,
133 we focus on the dataset obtained from the same region, which is more readily accessible.
134 Summing up, this study seeks to address the following hypotheses using multiple subsets of
135 the continuous United States (CONUS) dataset:

136

137 1. The availability of additional climate data, i.e. unlabeled data, could potentially enhance
138 the performance of LSTM models in producing reliable streamflow predictions in
139 diverse modeling scenarios. Therefore, implementing a semi-supervised learning-
140 based framework will be useful.



141 2. It would be beneficial to use a semi-supervised learning-based framework that
142 leverages both labeled and unlabeled data, but treats them differently instead of treating
143 them homogeneously. By differentiating between the two datasets and incorporating
144 them into the training process, the model can potentially achieve better performance
145 on unseen data.

146 3. The joint training of both labeled and unlabeled dataset has the potential to improve
147 model performance in comparison to a separate training approach (i.e., pre-training
148 followed by fine-tuning), which could also be employed to improve modeling
149 performance in data-sparse regions.

150
151 Through a series of experiments comparing our proposed semi-supervised learning-based
152 framework to diverse models, we aim to assess the hypotheses mentioned above and gain
153 insight into how LSTM models can enhance performance in data-sparse regions. This
154 exploration will enable us to better understand the benefits of our proposed framework.

155

156 **2. Methods and Data**

157 This section begins by introducing the dataset used in this study (section 2.1), followed by an
158 overview of the LSTM model structure (section 2.2) and the proposed framework based on the
159 semi-supervised learning (section 2.3). Finally, we outline the specific experimental designs
160 assessed in this study (section 2.4).

161

162 **2.1 Dataset**

163 To investigate the effectiveness of semi-supervised learning in analyzing the streamflow



164 network, which covers various geologies and climatic conditions, we use the Catchment
165 Attributes and MEteorology for Large-sample Studies (CAMELS) dataset (Newman et al.,
166 2015). We utilize the dataset due to their abundance of data, which could potentially strengthen
167 the validation of our hypotheses (addressed in section 2.4). The dataset includes the basin-
168 averaged hydrometeorological time series, catchment characteristics and daily streamflow
169 measurements for 671 basins over the CONUS. It is worth noting that the dataset has been
170 widely utilized to facilitate generalization and application of data-driven models for various
171 purposes (e.g., Feng et al., 2020; Gauch et al., 2021b; Kratzert et al., 2019b). We have adopted
172 the same subset of 531 basins as Gauch et al. (2021b) and Kratzert et al. (2019b) (see Figure 1)
173 while excluding 140 basins that display considerable inconsistencies in their calculated
174 watershed boundaries from different methodologies. Consistent with the aforementioned
175 studies, we utilize the Maurer meteorological forcing dataset, which includes daily cumulative
176 precipitation (*PRCP*), maximum and minimum air temperature (T_{max} and T_{min}), short-wave
177 radiation (*SRAD*), and vapor pressure (*VP*), spatially averaged for each basin. Furthermore, we
178 use 27 of the static catchment characteristics, including topography, climate characteristics,
179 land cover, soil and geology characteristics (Table 1). The spatially aggregated data have been
180 derived from an original gridded dataset that has a resolution of $1/8^\circ$. The meteorological
181 forcing and streamflow data are normalized so that all variables for each basin has a mean of
182 zero and unity variance.

183

184 **2.2 Long Short-Term Memory network**

185 In this work, we utilize a LSTM architecture for the rainfall-runoff modeling. A LSTM network
186 is a type of recurrent neural network designed to model long-term dependencies between input
187 and output data. LSTMs utilize an internal memory state that is updated at each time step by a



188 set of activated functions called *gates* (Hochreiter and Schmidhuber, 1997). The memory cells
189 are comparable to a state vector in a traditional dynamic system model, which leads to the
190 suitability of LSTMs for modeling dynamics such as rainfall-runoff relationships. Compared
191 to vanilla recurrent neural networks, LSTMs are less affected by the vanishing gradient issue
192 that has prevented effective model learning (Hochreiter and Schmidhuber, 1997). Given a raw
193 input sequence $\mathbf{x}_0 = [x_0^1, x_0^2, \dots, x_0^T]$ with T time steps, where each element \mathbf{x}_0^t is a vector
194 containing model input at time step t , we specifically employed the following equations for the
195 forward pass through the LSTM:

196

$$197 \quad \mathbf{x}^t = \mathfrak{D}(GELU(\mathbf{W}_x \mathbf{x}_0^t + \mathbf{b}_x)) \quad \text{Eq. (1)}$$

$$198 \quad \mathbf{i}^t = \sigma(\mathbf{W}_i \mathbf{x}^t + \mathbf{U}_i \mathbf{h}^{t-1}) + \mathbf{b}_i \quad \text{Eq. (2)}$$

$$199 \quad \mathbf{f}^t = \sigma(\mathbf{W}_f \mathbf{x}^t + \mathbf{U}_f \mathbf{h}^{t-1}) + \mathbf{b}_f \quad \text{Eq. (3)}$$

$$200 \quad \mathbf{g}^t = \tanh(\sigma(\mathbf{W}_g \mathbf{x}^t + \mathbf{U}_g \mathbf{h}^{t-1})) + \mathbf{b}_g \quad \text{Eq. (4)}$$

$$201 \quad \mathbf{o}^t = \sigma(\mathbf{W}_o \mathbf{x}^t + \mathbf{U}_o \mathbf{h}^{t-1}) + \mathbf{b}_o \quad \text{Eq. (5)}$$

$$202 \quad \mathbf{c}^t = \mathbf{g}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t \quad \text{Eq. (6)}$$

$$203 \quad \mathbf{h}^t = \tanh(\mathbf{c}^t) \odot \mathbf{o}^t \quad \text{Eq. (7)}$$

$$204 \quad \hat{\mathbf{y}}^t = \mathbf{W}_y \mathfrak{D}(\mathbf{h}^t) + \mathbf{b}_y \quad \text{Eq. (8)}$$

205

206 where \mathbf{i}^t , \mathbf{f}^t , \mathbf{o}^t , and \mathbf{g}^t are the input gate, forget gate, output gate, and cell input, respectively,
207 at time step t . The cell state and recurrent input are denoted by \mathbf{c}^t and \mathbf{h}^t . *GELU* refers to the
208 Gaussian Error Linear Units (Hendrycks and Gimpel, 2016). Also, two activation functions,
209 sigmoid and hyperbolic tangent, are denoted by σ and \tanh . \mathbf{W} , \mathbf{U} , and \mathbf{b} are learnable
210 parameters for each gate, where subscripts suggest which gate the weight vector is used for,



211 and \odot represents element-wise multiplication. The dropout operator is denoted by \mathfrak{D} , which
212 randomly sets some nodes along with corresponding of the network connections to zero in
213 training phase in order to reduce overfitting (Srivastava et al., 2014).

214

215 A linear embedding layer (Eq. 1) is utilized to preprocess the inputs before delivering them to
216 the LSTM cell, in order to prevent the possibility of critical inputs being dropped out by the \mathfrak{D}
217 operators and thereby reducing the model's performance.

218

219 **2.3 Semi-supervised learning-based framework to improve hydrologic prediction**

220 The proposed semi-supervised learning-based framework is an improved version of self-
221 training methods proposed in the ML field (Hinton et al., 2015; Yarowsky, 1995). Self-training
222 methods utilize unlabeled data by imputing predicted labels (called pseudo labels) to the
223 unlabeled data. Specifically, our approach is based on knowledge distillation, where the pseudo
224 labels for unlabeled dataset are generated from a pre-trained teacher model trained on labeled
225 dataset. Student model is trained in supervised manner on both the labeled and (pseudo label-
226 assigned) unlabeled datasets. In this work, both teacher and student models have the same
227 structure, as in self-distillation (Zhang et al., 2019).

228

229 Suppose that we are given a set of data \mathbb{D} including labeled data $\mathbb{L}^t = \{(\mathbf{x}_{0_1}^t, y_1^t), (\mathbf{x}_{0_2}^t, y_2^t),$
230 $\dots, (\mathbf{x}_{0_N}^t, y_N^t)\}$ at basin ($n = 1, \dots, N$) and time ($t = 1, \dots, T$) and unlabeled data $\mathbb{U}^t =$
231 $\{\mathbf{x}_{0_1}^t, \mathbf{x}_{0_2}^t, \dots, \mathbf{x}_{0_N}^t\}$ at time ($t = 1, \dots, T$). The framework requires two input datasets (\mathbb{L}^t and
232 \mathbb{U}^t). The labeled data \mathbb{L}^t is employed to train a teacher LSTM model by minimizing a loss
233 function. The teacher model is then used to estimate streamflow (i.e., pseudo streamflow \hat{y}_n^t)



234 on unlabeled data \mathbb{U}^t . Afterwards, we train a student LSTM model by minimizing the
235 combined loss on both labeled and unlabeled data with pseudo streamflow. Finally, we repeat
236 the process by reinstating the student as a teacher, which generates new pseudo streamflow,
237 allowing us to train a new student. A flowchart of the proposed semi-supervised learning-based
238 framework is presented in Figure 2. The algorithmic procedure for the framework is as follows.

239

240 [1] Train the teacher model $f^*(\cdot)$ by minimizing the loss function L on labeled data \mathbb{L}^t

$$241 \quad \frac{1}{N \times T} \sum_{n=1}^N \sum_{t=1}^T L(y_n^t, f^*(\mathbf{x}_{0_n}^t, \boldsymbol{\theta}^*)) \quad \text{Eq. (9)}$$

242

243 [2] Use the teacher model to generate pseudo streamflow for unlabeled data \mathbb{U}^t

$$244 \quad \hat{y}_n^t = f^*(\mathbf{x}_{0_n}^t, \boldsymbol{\theta}^*), \quad \forall t = 1, \dots, T \quad \text{Eq. (10)}$$

245

246 [3] Train the student model $f^{**}(\cdot)$ by minimizing the below loss function L to consider the
247 training balance between labeled and unlabeled data

$$248 \quad \frac{1}{N \times T} \sum_{n=1}^N \sum_{t=1}^T L(y_n^t, f^{**}(\mathbf{x}_{0_n}^t, \boldsymbol{\theta}^{**})) + \alpha(t) \frac{1}{N \times T} \sum_{n=1}^N \sum_{t=1}^T L(\hat{y}_n^t, f^{**}(\mathbf{x}_{0_n}^t, \boldsymbol{\theta}^{**})) \quad \text{Eq. (11)}$$

249

250 where $\alpha(t)$ is a balance coefficient at epoch t .

251

252 The suitability of $\alpha(t)$ significantly impacts the performance of the student model. When $\alpha(t)$
253 is high, the loss function is primarily influenced by \mathbb{U}^t , whereas a small value allows the
254 benefits from \mathbb{L}^t to become more apparent. Consequently, to mitigate the risk of ending up in
255 poor local optima, we employ an annealing process that incorporates a t -varying $\alpha(t)$ as
256 follows:



257

$$\alpha(t) = \begin{cases} \alpha_0 & t \leq J \\ 0 & J < t \end{cases} \quad \text{Eq. (12)}$$

259

260 where this study utilizes $\alpha_0 = 1$ and $J = 15$ based on the epoch adopted in this study.

261

262 [4] Further train the student model as a teacher model and go back to step [2]. Our experiment
263 involves two iterations where the student assumes the role of the teacher, but it may be
264 beneficial to conduct additional iterations.

265

266 2.4 Experiments

267 To depict the situation in hydrological data-sparse regions, this study considers two dimensions
268 for each research hypothesis. First, the Ψ subset of basins, rather than considering all 531
269 basins, are utilized under two differently defined regions (heterogeneous and homogeneous
270 regions). The approach is adopted since, in data-scarce regions, the numbers of the streamflow
271 gauge are also limited in reality. To explore the performance in heterogeneous regions, the Ψ
272 subset of basins over the CONUS are randomly selected, and the model performance is
273 investigated. The purpose of this analysis is to emulate the diverse environmental factors that
274 exist across the subspace of the considered area. This analysis is repeated 3 times to match the
275 experiment trial conducted in heterogeneous regions. For the analysis of homogeneous areas,
276 this study employs three regions, namely the North Atlantic, Southwest, and Southern Rockies
277 regions (Figure 1). The purpose of this analysis is to take into account a broad spectrum of
278 environmental conditions while reproducing a homogeneous environmental situation within
279 the target region. The North Atlantic region comprises 84 basins that are moderately affected



280 by snow accumulation and melting processes and are predominantly covered by dense forests
281 (with an average forest coverage of 89%). In contrast, the Southwest region comprises 66
282 basins with relatively flat topography and lesser snow influence compared to other regions.
283 Lastly, the Southern Rockies region (50 basins) is significantly influenced by snow process and
284 consists relatively arid catchments (aridity index 1.71 and average annual precipitation 700
285 mm/year). The three cases consider scenarios in which Ψ takes on values of 10, 30, and 50,
286 representing situations where the basin density network is relatively deficient, moderate, and
287 sufficient, respectively.

288

289 Second, we consider two training scenarios, single and multi-year training scenarios, to
290 represent the available data in data-scare regions. It is worth mentioning that in the data-scare
291 regions, streamflow records often have a restricted length. For the single-year training scenario,
292 all models are trained from October 1, 1988 to September 30, 1989 and validated from October
293 1, 1989 to September 30, 1990. For unlabeled extended data, we employ data from October 1,
294 1983 to September 30, 1988. The models are then evaluated over 12 years (October 1, 1996 to
295 September 30, 2008). For the multi-year training scenario, all models are trained for 3 years
296 (October 1, 1988 to September 30, 1991) and validated over 2 years (October 1, 1991 to
297 September 30, 1993). Afterwards, we use the same data in unlabeled and evaluation periods
298 adopted for the single-year training scenario.

299

300 **2.4.1 The effect of semi-supervised learning on individual and regional setting**

301 With the first experiment, this study evaluates our proposed framework if it bolsters the ability
302 of rainfall-runoff modeling. In particular, we hypothesize that implementing a semi-supervised
303 learning-based framework would yield benefits in a diverse model setting. To confirm this



304 hypothesis, we run this experiment with individual and regional model setting. For the
305 individual setting, we train one network separately for each basin (hereafter idv-LSTM). On
306 the other hand, for the regional setting, we train a regional scale single network using all data
307 across multiple basins while allowing the network to learn more general pattern of the input-
308 to-output relationship (hereafter rgn-LSTM). Although the benefits of the proposed framework
309 may be expected in idv-LSTM due to increased learning data, it is unclear whether there would
310 be additional benefits in a rgn-LSTM. Previous studies have shown that LSTM predictions are
311 reliable when the model is trained over a large set of basins and that regional models already
312 learn more general patterns from a diverse set of basins (e.g, Gauch et al., 2021b). Therefore,
313 it remains to be seen whether the proposed framework would offer additional benefits beyond
314 those already achieved by regional models trained on diverse basin data. In addition, we will
315 evaluate the performance of a regional model in this experiment by increasing the amount of
316 training data, specifically by including a larger number of basins in rgn-LSTM. This will allow
317 us to determine the maximum number of basins for which the proposed framework offers
318 additional benefits, once its effectiveness in a regional setting is confirmed. To establish a
319 comparison, we obtain simulation results from a LSTM using a standard train-validation-
320 testing framework. These results are then used as the baseline for evaluating the performance
321 of our proposed semi-supervised learning-based framework.

322

323 **2.4.2 The effect of the annealing process on the student model**

324 The objective of the second set of experiments is to examine the impact of the annealing process
325 (Eq. 12) adopted in the proposed framework. In simpler terms, we hypothesize that utilizing an
326 imbalanced-based cost function with both labeled and unlabeled data would enhance the
327 accuracy of the model. The rationale behind this is that by accounting for pseudo streamflow



328 and exploiting their impact, it enables the data-driven model to effectively learn the underlying
329 hydrologic response to input variables.

330

331 To investigate this, we employ variants of the rgn-LSTM model. Specifically, we additionally
332 develop five different versions of the rgn-LSTM model within a semi-supervised learning-
333 based framework. The first additional model (rgn-LSTM-vr1) treats the equivalent data for \mathbb{L}^t
334 and \mathbb{U}^t by replacing $\alpha(t)$ with a value of 1 and does not provide a distinguishing weight. The
335 remaining four models (rgn-LSTM-vr2, rgn-LSTM-vr3, rgn-LSTM-vr4, and rgn-LSTM-vr5)
336 use t -varying weight but adopt different formulations from Eq. 12. They are designed to
337 amplify, slowly increase, or slowly decrease the influence of \mathbb{U}^t . The specific $\alpha(t)$
338 configurations are presented in the supporting information (see Text S1). Also, Figure S3
339 shows how the t -varying weight is changed given increases to the epoch for each rgn-LSTM
340 model.

341

342 **2.4.3 Comparison of our proposed framework to the separate training approaches**

343 Previous studies have suggested the separate training approach as a means of improving neural
344 network models (Anderson and Radic, 2021; Read et al., 2019). In this approach, a model is
345 first pre-trained on a specific dataset to learn general patterns and relationships between input
346 and output data. The model is then fine-tuned on an additional dataset to learn more specific
347 behaviors and improve its performance on a particular task. This process allows the model to
348 adapt to the nuances of the task at hand, and has been shown to be effective in the ML field as
349 well (George et al., 2017; Yosinski et al., 2014).

350

351 For our third experiment, we aim to determine whether our proposed framework, which



352 incorporates joint training using both labeled and unlabeled data, can achieve better results
353 compared to the separate training approaches. To accomplish this, we introduce two additional
354 rgn-LSTM models: rgn-LSTM-sep and rgn-LSTM-trans. The rgn-LSTM-sep model initially
355 leverages unlabeled data U^t to capture the underlying patterns of runoff generating processes.
356 Subsequently, it undergoes fine-tuning on labeled data L^t to refine its performance, specifically
357 considering the delicate input and output relationships within specific basins. On the other hand,
358 the rgn-LSTM-trans model incorporates the recent technique proposed by Ma et al. (2021),
359 which utilizes Transfer Learning (Thrun and Pratt, 1998) (see Supplemental information for a
360 brief description of the method). They adopt a methodology wherein the models are initially
361 pretrained on a region abundant in data (known as the source region). These pretrained models
362 are subsequently transferred to data-scarce regions to overcome the limitations of local
363 observations. For this study, we employ the CAMELS-GB dataset, which is a comprehensive
364 dataset for Great Britain based on the CAMELS framework (Coxon et al., 2020), as our source
365 dataset. The dataset is selected because the CAMELS-GB basins exhibit a wide range of
366 hydrological conditions, analogous to the conditions found in our study basins. To be specific,
367 the rgn-LSTM-trans model is pretrained using 44 climate and basin attributes from the
368 CAMELS-GB dataset (as shown in Table S1).

369

370 **2.4.4 Evaluation metrics and hyperparameters**

371 To evaluate the modeling performance for each experiment, we run all models with four
372 random seeds and use the average estimated streamflow obtained from the resulting ensemble
373 members. The first metric used to assess the performance is the Nash-Sutcliffe efficiency (NSE)
374 coefficient (Nash and Sutcliffe, 1970), which is calculated for each basin. Also, we utilize two
375 metrics to evaluate the model's performance for both extreme flows: the modified Nash–



376 Sutcliffe efficiency (MNSE) and the logged transformed Nash-Sutcliffe efficiency (LNSE).
377 These metrics specifically focus on the performance of the model for high and low flows,
378 respectively (Ahn et al., 2016; Muleta, 2011). It is important to note that all of the metrics
379 reported in the manuscript are calculated based on the evaluation period.

380

381 Hyperparameters including learning rate, hidden states, length of input sequence, dropout rate,
382 epochs, and numbers of LSTM layer are configurations of LSTM model and thus yield varying
383 degrees of influence on the model's performance (Bengio et al., 2017). To avoid potential bias
384 in performance evaluation that may favor our proposed framework, we choose to adopt the
385 same hyperparameter configurations used in previous studies (Kratzert et al., 2021, 2019b),
386 rather than determining new ones for this study. Finally, all model configurations are trained
387 using the mean squared error (MSE) metric similar to the previous work.

388

389 **3. Results**

390 **3.1 Evaluating semi-supervised learning in data-scarce regions**

391 In this section, we assess the effectiveness of the proposed semi-supervised learning-based
392 framework in enhancing streamflow predictions. Figures 3 and 4 illustrate the spatial
393 distribution of the NSE difference for both idv-LSTM and rgn-LSTM cases, respectively,
394 during the evaluation period in comparison to the baseline models. Figures S2, S3 and S4
395 present the differences in other metrics (MNSE and LNSE) for the idv-LSTM and rgn-LSTM
396 settings. In each figure, the red color indicates that our proposed framework outperforms the
397 baseline models in terms of prediction accuracy, while the blue color indicates that our
398 proposed framework underperforms the baseline models. Additionally, Table 2 provides a
399 summary of the median performance across all experiments, encompassing the three evaluation



400 metrics. Notably, our framework exhibits improvements across all three metrics, underscoring
401 its effectiveness. The single-year training scenario in idv-LSTM stands out by yielding the
402 most significant benefits, with a notable improvement of 0.390 in median NSE and 0.238 in
403 MNSE. Similarly, the multi-year training scenario in idv-LSTM exhibits substantial
404 improvements, where our semi-supervised learning approach yields remarkable improvements
405 of 0.165 in median NSE and 0.374 in LNSE. The results indicate that our proposed framework
406 offers greater advantages when addressing regions with limited data availability, particularly
407 in data-scarce areas where the available data is relatively smaller.

408

409 Moreover, our proposed framework delivers substantial benefits in the context of rgn-LSTM.
410 When considering the single-year training scenario for the deficient network (i.e., $\Psi = 10$) in
411 rgn-LSTM across all six heterogeneous and homogeneous regions, it demonstrates a
412 remarkable improvement of 0.371 in median NSE, 0.275 in MNSE, and 0.560 in LNSE.
413 Similarly, in the case of the sufficient network (i.e., $\Psi = 50$), the multi-year training scenario
414 yields an improvement of 0.023 in median NSE, 0.027 in MNSE, and 0.062 in LNSE. The
415 results reveal several notable insights. Firstly, similar to idv-LSTM, our framework
416 demonstrates increased effectiveness when dealing with insufficient records. This highlights
417 its utility in situations where data availability is significantly limited. We also note that, for
418 some basins particularly in the Southern Rockies region, the baseline model performs better
419 than the models trained by our framework (see Figure 4). The performance declines may be
420 related to frequently having zero discharge in observation. Having zero values for a high
421 percentage of the training samples seems to be a difficult information for the teacher model to
422 learn and to reproduce this hydrological behavior and affect the performance of the student
423 model. However, we observe that the median of a metric is still positive, indicating that the



424 models trained by our framework performs are effective. Next, the efficacy of our framework
425 extends to relatively large streamflow networks, as evidenced by our results in a network
426 comprising 50 basins. Our proposed framework offers additional benefits that surpass those
427 achieved by regional models trained on diverse basin data, even when rgn-LSTM has already
428 learned general patterns from a diverse set of basins. This is particularly relevant in data-scarce
429 regions where some streamflow stations may be available with limited records. Consequently,
430 the findings from this analysis provide valuable insights that can guide the practical
431 implementation of our framework in real-world applications, addressing the challenges posed
432 by data scarcity in streamflow prediction.

433

434 It is important to highlight that the effectiveness of our proposed framework is especially
435 pronounced when using a separate network for each basin (idv-LSTM). Also, there is an
436 expectation that rgn-LSTM would still exhibit improvement when utilizing a semi-supervised
437 learning-based framework. This suggests that employing a single setting for our remaining
438 assessment is acceptable. Furthermore, as previously mentioned, it is probable that some
439 streamflow stations are available even in data-scarce regions. This suggests that conducting an
440 analysis by combining data from those stations with regional models trained on multiple basin
441 data would offer a more realistic evaluation. Therefore, for the remaining analysis, we will
442 adopt rgn-LSTM particularly with the moderate density network.

443

444 **3.2 Evaluating the selection of the annealing process on the student model**

445 In the first experiment, we confirm the benefits of a semi-supervised learning-based framework
446 in enhancing streamflow predictions. We now analyze the performance of six rgn-LSTM
447 models (see Figure S1) to explore the appropriateness of the annealing process (Eq. 12) adopted



448 in the proposed framework. Figures 5 and S5 show that results of rgn-LSTM obtained over the
449 evaluation period with the single and multi-year training scenarios, respectively.

450

451 Those figures present two noteworthy observations. First, the performance of the rgn-LSTM-
452 vr1 model is notably lower compared to the other models. Specifically, significant declines in
453 performance are observed for the single-year training scenario, while its performance remains
454 similar to the other models for the multi-year training scenario. However, even in the multi-
455 year training scenario, lower performance is evident in LNSE particularly for the
456 heterogeneous region. These findings suggest that incorporating an imbalanced-based cost
457 function between labeled and unlabeled data enhances the model's predictive capabilities. Next,
458 the models employing structures that diminish the influence of unlabeled data (e.g., rgn-LSTM-
459 vr5 and rgn-LSTM) show better results compared to the models that amplify the role of
460 unlabeled data (e.g., rgn-LSTM-vr2 and rgn-LSTM-vr4) particularly in the single-year training
461 scenarios. There could be multiple factors contributing to this disparity, but our inference is
462 that the outperformance may be attributed to the low quality in learning of the teacher model
463 due to insufficient data. The low quality for the teacher model potentially affects the quality of
464 the unlabeled data. By leveraging the expanded training data that includes unlabeled data, the
465 student model can gain a rough understanding of streamflow modeling. This initial exploration
466 proves beneficial, allowing the model to converge quickly and reducing the chances of
467 overfitting. Subsequently, the network undergoes fine-tuning using high-quality labeled data
468 on the latter part of the epoch progression. Therefore, the models employing structures to
469 diminish the influence of unlabeled data would be beneficial. Our inference is also supported
470 by the multi-year training scenario. While rgn-LSTM remains competitive, its superiority
471 becomes less apparent due to the improved learning of the teacher model resulting from the



472 expanded training samples. Summing up, by differentiating between the two unlabeled and
473 labeled datasets, the model can potentially achieve better performance in a semi-supervised
474 learning-based framework.

475

476 **3.3 Comparing our proposed model with rgn-LSTM-sep and rgn-LSTM-trans**

477 Finally, we compare our semi-supervised learning-based framework with two separate training
478 approaches, namely rgn-LSTM-sep and rgn-LSTM-trans. Figures 6 and S6 present the spatial
479 distribution of the NSE, MNSE, and LNSE metric differences for the three regions in the
480 heterogeneous and homogeneous regions, respectively. The figures illustrate the relative
481 performance of our framework's models compared to the two fine-tuning models. Here, the
482 utilization of the red color highlights instances where our proposed framework surpasses a
483 separate training approach in terms of prediction accuracy in the evaluation period.

484

485 Based on the comparison between rgn-LSTM and rgn-LSTM-sep, the benefits of utilizing a
486 semi-supervised learning approach over relying solely on weight initialization using unlabeled
487 data (corresponding to pre-training) are evident. For example, the single-year training scenario
488 in heterogeneous region yields notable benefits, with an improvement of 0.024 in median NSE
489 and 0.031 in MNSE. Similarly, the multi-year training scenario also show substantial
490 improvements, where our semi-supervised learning approach yields remarkable improvements
491 of 0.020 in median NSE and 0.044 in MNSE. Particularly, the most significant improvement
492 is observed in mean LNSEs for both scenarios, with our proposed framework achieving a
493 noteworthy improvement of 0.077 and 0.090 when compared to rgn-LSTM-sep. The results
494 indicate that the joint training of both labeled and unlabeled datasets leads to better
495 performance than the separate training approach that utilizes weight initialization only with



496 unlabeled data.

497

498 When comparing rgn-LSTM and rgn-LSTM-trans, we observe slight differences in contrast to
499 the results between rgn-LSTM and rgn-LSTM-sep, suggesting that those approaches (rgn-
500 LSTM and rgn-LSTM-trans) provide reliable predictions in the rainfall-runoff relationship.
501 However, rgn-LSTM tends to exhibit higher prediction accuracy overall compared to rgn-
502 LSTM-trans. This is especially evident when considering the MNSE and LNSE metrics,
503 highlighting the effectiveness of a semi-supervised learning-based framework in more
504 accurately representing local extreme flows. One possible explanation for the performance
505 difference between rgn-LSTM and rgn-LSTM-trans is that transfer learning can be effective
506 when the dataset in both the source and target regions are sufficiently similar. However, rgn-
507 LSTM-trans is based on pre-trained knowledge from the source region, utilizing 44 forcing
508 variables that substantially differ from the attributes employed in this study (refer to Table 1).
509 It is commonly referred to as negative transfer learning (Torrey and Shavlik, 2010). Our
510 inference is supported by the outperformance of rgn-LSTM in the Southwest region, which
511 exhibits conditions fairly different from those of Great Britain (see Figure S6). It is also worth
512 noting in this comparison that rgn-LSTM may have a disadvantage due to the fact that rgn-
513 LSTM-trans is trained using a 10-year labeled dataset, which is nearly double the additional
514 data used in the training process for rgn-LSTM. Taken together, we therefore consider a semi-
515 supervised learning-based approach a useful and complementary approach to the transfer
516 learning approach, but would caution against using it as a replacement for bolstering the ability
517 of rainfall-runoff modeling in all cases.

518

519 **4. Discussion**



520 **4.1 Impact of the performance of teacher model**

521 The predictive capability of the teacher model is vital within the proposed semi-supervised
522 learning-based framework. This is because the teacher model is employed to generate pseudo
523 streamflow on unlabeled data. Consequently, an enhanced performance of the teacher model is
524 anticipated to result in greater improvements within the framework. Figure 7 shows the
525 accuracy improvements obtained in the semi-supervised learning-based framework relative to
526 the baselines when they are compared to the performance of the teacher model. The figure
527 presents results for two model settings, idv-LSTM and rgn-LSTM with the moderate density
528 network. It is important to note that similar patterns are observed in the other results (not
529 shown). Interestingly, the impact of the performance of the teacher model is different from our
530 expectation. While there are slight variations in each plot, the anticipated improvement
531 (indicated by the red lines) generally follows an upward trend, reaching its peak around a NSE
532 value of 0.4, and subsequently experiencing a decline in improvement.

533

534 These findings suggest that achieving higher performance in the teacher model does not
535 necessarily translate into greater improvement within the framework. One potential
536 explanation is that the involvement of numerous latent processes in the rainfall-runoff process.
537 These processes include factors such as subsurface interactions (e.g., aquifer dynamics and
538 transmissivity). Due to the complexity of these confounding factors, it becomes challenging
539 for the network to further capture the entire runoff generation process especially in the basins
540 well trained by the student network. Instead, our analysis shows that the proposed framework
541 exhibits its highest effectiveness in the study basins when the network achieves a moderate
542 level of accuracy, specifically around an NSE value of approximately 0.4 when the baseline
543 network is applied.



544

545 **4.2 Applicability of the semi-supervised learning-based framework**

546 While semi-supervised learning holds significant promise, further exploration is encouraged to
547 assess its applicability, particularly in the context of time series tasks. The successful
548 application of semi-supervised learning has predominantly been observed in computer vision
549 tasks (Chen et al., 2021; Yang et al., 2019), but it has also demonstrated success in other
550 machine learning domains (Wang et al., 2021; Zhu et al., 2021). However, the success has
551 rarely been extended to time series related tasks (Dai et al., 2023). This scarcity of success in
552 time series tasks further underscores the significance and value of the present study.

553

554 Furthermore, the improvement of hydrological modeling initiatives has been dependent on both
555 sufficient data collection and enhancements in the model's algorithm to an equal degree. The
556 sparsity and inconsistency of the meteorological dataset additionally result in low performance
557 in the streamflow prediction and create a problematic situation to implement our proposed
558 framework. In developing countries, the situation arises due to the insufficient availability of
559 equipment used for monitoring meteorological data, such as precipitation and air temperature.
560 To be specific, the lack of sufficient data for tracking meteorological information in African
561 countries contributes to the encountered situation. Although initiatives like The trans-African
562 hydrometeorological observatory (TAHMO) have been launched (van de Giesen et al., 2014),
563 there is still a significant gap in data availability and coverage. As a potential solution to
564 mitigate the situation, we can consider the utilization of reanalysis-based climate data, such as
565 the global dataset provided by the European Centre for Medium-Range Weather Forecasts
566 (ECMWF). Additionally, employing approaches like statistical downscaling of these global



567 datasets, as demonstrated in the studies by Voropay et al. (2021) and Xie et al. (2022), could
568 prove effective in mitigating the challenges presented by the limited meteorological data.

569

570 **4.3 Future work**

571 This paper utilizes a semi-supervised learning approach to improve the predictive ability in
572 rainfall-runoff modeling while addressing the limitations associated with streamflow
573 observations. The proposed framework utilizes predicted streamflow estimated by a teacher
574 model as pseudo-labels, indicating high-quality pseudo-labels is important for the performance
575 of the student model. However, this study does not address the issue of uncertainty associated
576 with these pseudo-labels. One potential solution is to employ Bayesian neural networks (BNNs;
577 Kendall and Gal, 2017), which effectively handle input data noise, known as aleatoric
578 uncertainty, by incorporating its impact into the loss function. This utilization of BNNs as a
579 heteroscedastic modeling technique may be useful to reduce prediction variance and enhance
580 the quality of pseudo-labels obtained. Our team intends to explore this approach in the near
581 future as part of our ongoing research efforts.

582

583 **5. Conclusions**

584 The science of hydrology has primarily evolved by leveraging established physical and
585 empirical relationships to comprehend the complex dynamics of rainfall-runoff interactions.
586 Although significant progress has been made in harnessing data-driven models to enhance
587 insights and intuition derived from abundant hydrological dataset, a fundamental obstacle
588 remains due to the scarcity of available data.

589



590 In this study, we developed a semi-supervised learning-based framework to mitigate the
591 challenges associated with predicting streamflow in regions with limited data availability. The
592 framework enables a data-driven model to enhance its training dataset by incorporating
593 additional climate data, even in scenarios with limited paired records of climate and streamflow
594 data. This is achieved through the generation of pseudo streamflow data. In particular, we
595 introduced a novel loss function for the student model, designed to effectively distinguish the
596 contributions of labeled and unlabeled data to the loss function during the training process.
597 Through a range of diverse experimental designs, we conducted extensive validation to
598 demonstrate the substantial efficacy of the proposed framework in comparison to a simple
599 baseline model. Lastly, we conducted a thorough comparison between our proposed framework
600 and two separate training approaches, affirming the effectiveness of our framework. We firmly
601 believe that the value of this framework is immense, as it capitalizes on the availability of
602 longer historical climate data records, including the utilization of global climate datasets. This
603 is particularly advantageous in regions where streamflow records are scarce, as it facilitates the
604 extraction of valuable insights from the wealth of accessible climate data.

605

606 **Code availability**

607 The code is available upon the request to the corresponding author.

608

609 **Author contribution**

610 Sunghyun Yoon carried out data acquisition, formal analysis, investigation, visualization, and
611 writing. Kuk-Hyun Ahn conceptualized the study, designed the methodology, performed



612 formal analysis, participated in writing the draft, and served as supervisor of the study.

613

614 **Competing interests**

615 The authors declare that they have no conflict of interest.

616

617 **Acknowledgements**

618 This work was supported by the National Research Foundation of Korea(NRF) grant funded
619 by the Korea government(MSIT) (No. RS-2023-00208210).

620

621 **References**

- 622 Abbott, M.B., 1999. Introducing hydroinformatics. *J. Hydroinformatics* 1, 3–19.
- 623 Ahn, K.-H., 2021. Streamflow estimation at partially gaged sites using multiple-dependence
624 conditions via vine copulas. *Hydrol. Earth Syst. Sci.* 25, 4319–4333.
- 625 Ahn, K.-H., Steinschneider, S., Palmer, R., 2016. A hierarchical Bayesian model for
626 regionalized seasonal forecasts of low flows in the northeastern United States. *Water*
627 *Resour. Res.*
- 628 Alquraish, M.M., Khadr, M., 2021. Remote-Sensing-Based Streamflow Forecasting Using
629 Artificial Neural Network and Support Vector Machine Models. *Remote Sens.* 13, 4147.
- 630 Anderson, S., Radic, V., 2021. Evaluation and interpretation of convolutional-recurrent
631 networks for regional hydrological modelling. *Hydrol Earth Syst Sci Discusspreprint*
632 [Httsdoi Org105194hess-2021-113](https://doi.org/10.5194/hess-2021-113) Rev.
- 633 Arsenault, R., Martel, J.-L., Brunet, F., Brissette, F., Mai, J., 2023. Continuous streamflow
634 prediction in ungauged basins: long short-term memory neural networks clearly
635 outperform traditional hydrological models. *Hydrol. Earth Syst. Sci.* 27, 139–157.
- 636 Bartoletti, N., Casagli, F., Marsili-Libelli, S., Nardi, A., Palandri, L., 2018. Data-driven
637 rainfall/runoff modelling based on a neuro-fuzzy inference system. *Environ. Model.*
638 *Softw.* 106, 35–47.
- 639 Bengio, Y., Goodfellow, I., Courville, A., 2017. *Deep learning*. MIT press Cambridge, MA,
640 USA.
- 641 Beven, K., 2020. Deep learning, hydrological processes and the uniqueness of place. *Hydrol.*
642 *Process.* 34, 3608–3613.
- 643 Beven, K.J., 2011. *Rainfall-runoff modelling: the primer*. John Wiley & Sons.



- 644 Birkinshaw, S.J., James, P., Ewen, J., 2010. Graphical user interface for rapid set-up of
645 SHETRAN physically-based river catchment model. *Environ. Model. Softw.* 25, 609–
646 610.
- 647 Bitew, M.M., Gebremichael, M., 2011. Assessment of satellite rainfall products for streamflow
648 simulation in medium watersheds of the Ethiopian highlands. *Hydrol. Earth Syst. Sci.*
649 15, 1147–1155.
- 650 Booker, D., Woods, R., 2014. Comparing and combining physically-based and empirically-
651 based approaches for estimating the hydrology of ungauged catchments. *J. Hydrol.* 508,
652 227–239.
- 653 Bowes, B.D., Sadler, J.M., Morsy, M.M., Behl, M., Goodall, J.L., 2019. Forecasting
654 groundwater table in a flood prone coastal city with long short-term memory and
655 recurrent neural networks. *Water* 11, 1098.
- 656 Burnash, R.J., Ferral, R.L., McGuire, R.A., 1973. A generalized streamflow simulation system,
657 conceptual modeling for digital computers.
- 658 Chadalawada, J., Herath, H., Babovic, V., 2020. Hydrologically informed machine learning for
659 rainfall-runoff modeling: A genetic programming-based toolkit for automatic model
660 induction. *Water Resour. Res.* 56, e2019WR026933.
- 661 Chen, I.-T., Chang, L.-C., Chang, F.-J., 2018. Exploring the spatio-temporal interrelation
662 between groundwater and surface water by using the self-organizing maps. *J. Hydrol.*
663 556, 131–142.
- 664 Chen, X., Yuan, Y., Zeng, G., Wang, J., 2021. Semi-supervised semantic segmentation with
665 cross pseudo supervision, in: *Proceedings of the IEEE/CVF Conference on Computer
666 Vision and Pattern Recognition.* pp. 2613–2622.
- 667 Chen, Y., Mancini, M., Zhu, X., Akata, Z., 2022. Semi-supervised and unsupervised deep visual
668 learning: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*
- 669 Coxon, G., Addor, N., Bloomfield, J., Freer, J., Fry, M., Hannaford, J., Howden, N., Lane, R.,
670 Lewis, M., Robinson, E., others, 2020. Catchment attributes and hydro-meteorological
671 timeseries for 671 catchments across Great Britain (CAMELS-GB).
- 672 Dai, W., Li, X., Cheng, K.-T., 2023. Semi-Supervised Deep Regression with Uncertainty
673 Consistency and Variational Model Ensembling via Bayesian Neural Networks. *ArXiv
674 Prepr. ArXiv230207579.*
- 675 Do, H.X., Westra, S., Leonard, M., 2017. A global-scale investigation of trends in annual
676 maximum streamflow. *J. Hydrol.* 552, 28–43.
- 677 Du, F., Zhu, A.-X., Liu, J., Yang, L., 2020. Predictive mapping with small field sample data
678 using semi-supervised machine learning. *Trans. GIS* 24, 315–331.
- 679 Feng, D., Fang, K., Shen, C., 2020. Enhancing streamflow forecast and extracting insights
680 using long-short term memory networks with data integration at continental scales.
681 *Water Resour. Res.* 56, e2019WR026793.
- 682 Frame, J.M., Kratzert, F., Raney, A., Rahman, M., Salas, F.R., Nearing, G.S., 2021. Post-
683 processing the national water model with long short-term memory networks for
684 streamflow predictions and model diagnostics. *JAWRA J. Am. Water Resour. Assoc.*
685 57, 885–905.
- 686 Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., Hochreiter, S., 2021a. Rainfall–runoff
687 prediction at multiple timescales with a single Long Short-Term Memory network.
688 *Hydrol. Earth Syst. Sci.* 25, 2045–2062.
- 689 Gauch, M., Mai, J., Lin, J., 2021b. The proper care and feeding of CAMELS: How limited
690 training data affects streamflow prediction. *Environ. Model. Softw.* 135, 104926.



- 691 George, D., Shen, H., Huerta, E., 2017. Deep Transfer Learning: A new deep learning glitch
692 classification method for advanced LIGO. ArXiv Prepr. ArXiv170607446.
- 693 Han, H., Choi, C., Jung, J., Kim, H.S., 2021. Deep learning with long short term memory based
694 sequence-to-sequence model for rainfall-runoff simulation. *Water* 13, 437.
- 695 He, K., Fan, H., Wu, Y., Xie, S., Girshick, R., 2020. Momentum contrast for unsupervised
696 visual representation learning, in: Proceedings of the IEEE/CVF Conference on
697 Computer Vision and Pattern Recognition. pp. 9729–9738.
- 698 Hendrycks, D., Gimpel, K., 2016. Gaussian error linear units (gelus). ArXiv Prepr.
699 ArXiv160608415.
- 700 Hinton, G., Vinyals, O., Dean, J., 2015. Distilling the knowledge in a neural network. ArXiv
701 Prepr. ArXiv150302531.
- 702 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Comput.* 9, 1735–1780.
- 703 Hoedt, P.-J., Kratzert, F., Klotz, D., Halmich, C., Holzleitner, M., Nearing, G.S., Hochreiter, S.,
704 Klambauer, G., 2021. Mc-lstm: Mass-conserving lstm, in: International Conference on
705 Machine Learning. PMLR, pp. 4275–4286.
- 706 Hunt, K.M., Matthews, G.R., Pappenberger, F., Prudhomme, C., 2022. Using a long short-term
707 memory (LSTM) neural network to boost river streamflow forecasts over the western
708 United States. *Hydrol. Earth Syst. Sci.* 26, 5449–5472.
- 709 Kendall, A., Gal, Y., 2017. What uncertainties do we need in bayesian deep learning for
710 computer vision? *Adv. Neural Inf. Process. Syst.* 30.
- 711 Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G.,
712 Hochreiter, S., Nearing, G., 2022. Uncertainty estimation with deep learning for
713 rainfall–runoff modeling. *Hydrol. Earth Syst. Sci.* 26, 1673–1693.
- 714 Kothari, M., Gharde, K., 2015. Application of ANN and fuzzy logic algorithms for streamflow
715 modelling of Savitri catchment. *J. Earth Syst. Sci.* 124, 933–943.
- 716 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., Herrnegger, M., 2018. Rainfall–runoff
717 modelling using long short-term memory (LSTM) networks. *Hydrol. Earth Syst. Sci.*
718 22, 6005–6022.
- 719 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A.K., Hochreiter, S., Nearing, G.S., 2019a.
720 Toward improved predictions in ungauged basins: Exploiting the power of machine
721 learning. *Water Resour. Res.* 55, 11344–11354.
- 722 Kratzert, F., Klotz, D., Hochreiter, S., Nearing, G.S., 2021. A note on leveraging synergy in
723 multiple meteorological data sets with deep learning for rainfall–runoff modeling.
724 *Hydrol. Earth Syst. Sci.* 25, 2685–2703.
- 725 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., Nearing, G., 2019b. Towards
726 learning universal, regional, and local hydrological behaviors via machine learning
727 applied to large-sample datasets. *Hydrol. Earth Syst. Sci.* 23, 5089–5110.
- 728 Lee, D.-G., Ahn, K.-H., 2022. Assessment of Suitable Gridded Climate Datasets for Large-
729 Scale Hydrological Modelling over South Korea. *Remote Sens.* 14, 3535.
- 730 Lees, T., Buechel, M., Anderson, B., Slater, L., Reece, S., Coxon, G., Dadson, S.J., 2021.
731 Benchmarking Data-Driven Rainfall-Runoff Models in Great Britain: A comparison of
732 LSTM-based models with four lumped conceptual models. *Hydrol. Earth Syst. Sci.*
- 733 Leisher, C., Makau, J., Kihara, F., Kariuki, A., Sowles, J., Courtemanch, D., Njugi, G., Apse,
734 C., 2016. Upper Tana-Nairobi Water Fund Monitoring and Evaluation Plan. Nat.
735 Conserv. IFAD CIAT GEF TRICOKEN.
- 736 Levatić, J., Ceci, M., Kocев, D., Džeroski, S., 2017. Self-training for multi-target regression
737 with tree ensembles. *Knowl.-Based Syst.* 123, 41–60.



- 738 Ley, A., Bormann, H., Casper, M., 2023. Intercomparing LSTM and RNN to a Conceptual
739 Hydrological Model for a Low-Land River with a Focus on the Flow Duration Curve.
740 *Water* 15, 505.
- 741 Li, D., Marshall, L., Liang, Z., Sharma, A., Zhou, Y., 2021. Bayesian LSTM with stochastic
742 variational inference for estimating model uncertainty in process-based hydrological
743 models. *Water Resour. Res.* 57, e2021WR029772.
- 744 Liang, C., Li, H., Lei, M., Du, Q., 2018. Dongting lake water level forecast and its relationship
745 with the three gorges dam based on a long short-term memory network. *Water* 10, 1389.
- 746 Liang, X., Lettenmaier, D.P., Wood, E.F., Burges, S.J., 1994. A simple hydrologically based
747 model of land surface water and energy fluxes for general circulation models. *J.*
748 *Geophys. Res. Atmospheres* 99, 14415–14428.
- 749 Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., Sharma, A., Shen, C., 2021.
750 Transferring hydrologic data across continents—leveraging data-rich regions to improve
751 hydrologic prediction in data-sparse regions. *Water Resour. Res.* 57, e2020WR028600.
- 752 Mckane, R., Brookes, A., Djang, K., Stieglitz, M., Abdelnour, A., Halama, J., Pettus, P., Phillips,
753 D., 2014. VELMA Version 2.0 User Manual and Technical Documentation. Corvallis
754 Or. https://www.epa.gov/sites/production/files/2016-01/documents/velma2_0usermanual.pdf
755 Pdflast Accessed 1006 17.
- 756 Muleta, M.K., 2011. Model performance sensitivity to objective function during automated
757 calibrations. *J. Hydrol. Eng.* 17, 756–767.
- 758 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I — A
759 discussion of principles. *J. Hydrol.* 10, 282–290. [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)
- 761 Nearing, G.S., Kratzert, F., Sampson, A.K., Pelissier, C.S., Klotz, D., Frame, J.M., Prieto, C.,
762 Gupta, H.V., 2021. What role does hydrological science play in the age of machine
763 learning? *Water Resour. Res.* 57, e2020WR028091.
- 764 Newman, A., Clark, M., Sampson, K., Wood, A., Hay, L., Bock, A., Viger, R., Blodgett, D.,
765 Brekke, L., Arnold, J., others, 2015. Development of a large-sample watershed-scale
766 hydrometeorological data set for the contiguous USA: data set characteristics and
767 assessment of regional variability in hydrologic model performance. *Hydrol. Earth Syst.*
768 *Sci.* 19, 209–223.
- 769 Nourani, V., Komasi, M., Mano, A., 2009. A multivariate ANN-wavelet approach for rainfall–
770 runoff modeling. *Water Resour. Manag.* 23, 2877–2894.
- 771 Oruche, R., Egede, L., Baker, T., O’Donncha, F., 2021. Transfer learning to improve
772 streamflow forecasts in data sparse regions. *ArXiv Prepr. ArXiv211203088*.
- 773 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., 2019. Deep
774 learning and process understanding for data-driven Earth system science. *Nature* 566,
775 195–204.
- 776 Roy, B., Singh, M.P., Kaloop, M.R., Kumar, D., Hu, J.-W., Kumar, R., Hwang, W.-S., 2021.
777 Data-Driven Approach for Rainfall-Runoff Modelling Using Equilibrium Optimizer
778 Coupled Extreme Learning Machine and Deep Neural Network. *Appl. Sci.* 11, 6238.
- 779 Seibert, J., Vis, M.J., 2012. Teaching hydrological modeling with a user-friendly catchment-
780 runoff-model software package. *Hydrol. Earth Syst. Sci.* 16, 3315–3325.
- 781 Shen, C., 2018. A transdisciplinary review of deep learning research and its relevance for water
782 resources scientists. *Water Resour. Res.* 54, 8558–8593.
- 783 Shen, C., Chen, X., Laloy, E., 2021. Broadening the use of machine learning in hydrology.
784 *Front. Water.*



- 785 Sitterson, J., Knightes, C., Parmar, R., Wolfe, K., Avant, B., Muche, M., 2018. An overview of
786 rainfall-runoff model types.
- 787 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a
788 simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15, 1929–
789 1958.
- 790 Taormina, R., Chau, K.-W., 2015. Data-driven input variable selection for rainfall–runoff
791 modeling using binary-coded particle swarm optimization and Extreme Learning
792 Machines. *J. Hydrol.* 529, 1617–1632.
- 793 Thrun, S., Pratt, L., 1998. Learning to learn: Introduction and overview. *Learn. Learn* 3–17.
- 794 TNC, 2015. Upper Tana-Nairobi water fund business case.
- 795 Torrey, L., Shavlik, J., 2010. Transfer learning, in: *Handbook of Research on Machine Learning*
796 *Applications and Trends: Algorithms, Methods, and Techniques.* IGI global, pp. 242–
797 264.
- 798 van de Giesen, N., Hut, R., Selker, J., 2014. The trans-African hydro-meteorological
799 observatory (TAHMO). *Wiley Interdiscip. Rev. Water* 1, 341–348.
- 800 Van, S.P., Le, H.M., Thanh, D.V., Dang, T.D., Loc, H.H., Anh, D.T., 2020. Deep learning
801 convolutional neural network in rainfall–runoff modelling. *J. Hydroinformatics* 22,
802 541–561.
- 803 Voropay, N., Ryazanova, A., Dyukarev, E., 2021. High-resolution bias-corrected precipitation
804 data over South Siberia, Russia. *Atmospheric Res.* 254, 105528.
- 805 Wang, C., Riviere, M., Lee, A., Wu, A., Talnikar, C., Haziza, D., Williamson, M., Pino, J.,
806 Dupoux, E., 2021. Voxpopuli: A large-scale multilingual speech corpus for
807 representation learning, semi-supervised learning and interpretation. *ArXiv Prepr.*
808 *ArXiv210100390.*
- 809 Wi, S., Steinschneider, S., 2022. Assessing the physical realism of deep learning hydrologic
810 model projections under climate change. *Water Resour. Res.* 58, e2022WR032123.
- 811 Xiang, Z., Yan, J., Demir, I., 2020. A rainfall-runoff model with LSTM-based sequence-to-
812 sequence learning. *Water Resour. Res.* 56, e2019WR025326.
- 813 Xie, K., Liu, P., Zhang, J., Han, D., Wang, G., Shen, C., 2021. Physics-guided deep learning
814 for rainfall-runoff modeling by considering extreme events and monotonic relationships.
815 *J. Hydrol.* 603, 127043.
- 816 Xie, W., Yi, S., Leng, C., Xia, D., Li, M., Zhong, Z., Ye, J., 2022. The evaluation of IMERG
817 and ERA5-Land daily precipitation over China with considering the influence of gauge
818 data bias. *Sci. Rep.* 12, 8085.
- 819 Xu, Y., Hu, C., Wu, Q., Jian, S., Li, Z., Chen, Y., Zhang, G., Zhang, Z., Wang, S., 2022.
820 Research on particle swarm optimization in LSTM neural networks for rainfall-runoff
821 simulation. *J. Hydrol.* 608, 127553.
- 822 Yang, T.-Y., Chen, Y.-T., Lin, Y.-Y., Chuang, Y.-Y., 2019. Fsa-net: Learning fine-grained
823 structure aggregation for head pose estimation from a single image, in: *Proceedings of*
824 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* pp. 1087–
825 1096.
- 826 Yarowsky, D., 1995. Unsupervised word sense disambiguation rivaling supervised methods, in:
827 *33rd Annual Meeting of the Association for Computational Linguistics.* pp. 189–196.
- 828 Yosinski, J., Clune, J., Bengio, Y., Lipson, H., 2014. How transferable are features in deep
829 neural networks? *Adv. Neural Inf. Process. Syst.* 27.
- 830 Zhang, L., Song, J., Gao, A., Chen, J., Bao, C., Ma, K., 2019. Be your own teacher: Improve
831 the performance of convolutional neural networks via self distillation, in: *Proceedings*



832 of the IEEE/CVF International Conference on Computer Vision. pp. 3713–3722.
833 Zhou, Z.-H., Zhou, Z.-H., 2021. Semi-supervised learning. *Mach. Learn.* 315–341.
834 Zhu, Y., Shareghi, E., Li, Y., Reichart, R., Korhonen, A., 2021. Combining Deep Generative
835 Models and Multi-lingual Pretraining for Semi-supervised Document Classification.
836 ArXiv Prepr. ArXiv210110717.

837 **List of Figures**

838 Figure 1 Locations of the selected 531 gauges in the United States. The gauge stations
839 across three homogenous regions are also illustrated.

840

841 Figure 2 Overview of the semi-supervised learning-based framework proposed in this
842 study. Here, the colored geometric shapes represent the recorded samples, pairing climate
843 forcing data with streamflow.

844

845 Figure 3 Difference of NSE results of idv-LSTM compared to their baseline models for (a)
846 single and (b) multi-year training scenarios. The color maps are limited to [-0.2, 1.0] for
847 the single training scenario and [-0.1, 0.5] for the multi-year training scenario for enhanced
848 visualization.

849

850 Figure 4 Difference of NSE results of rgn-LSTM compared to their baseline models across
851 experimental factors including three defined regions, two training scenarios, and three
852 basin densities in network. Here, the median NSE differences across basins in three defined
853 regions are presented in each plot.

854

855 Figure 5 Cumulative density functions of the results of the annealing process on rgn-LSTM
856 with the single-year training scenarios obtained for basins across (a), (b), (c) heterogeneous
857 region; and (d), (e), (f) homogeneous region. Here, three metrics, namely NSE (first
858 column), MNSE (second column), and LNSE (last column), are utilized.

859

860 Figure 6 Difference of performance in the three metrics, NSE (first column), MNSE
861 (second column), and LNSE (third column), of rgn-LSTM compared to the two fine-tuning
862 approaches (rgn-LSTM-sep and rgn-LSTM-trans) across three basin networks in
863 heterogeneous regions. Here, the median NSE differences across basins in three defined
864 regions are presented in each plot.

865

866 Figure 7 Plots comparing the NSEs achieved by the teacher model during the validation



867 period and the improved NSEs between the student models and baseline models during
868 the evaluation period under two model setting: (a) idv-LSTM and (b) rgn-LSTM with the
869 moderate density network, in multi-year training scenarios. Additionally, the expected
870 improved NSEs corresponding to the NSE of the teacher model are depicted, along with a
871 Lowess fit represented by a red line.

872

873 **List of Tables**

874

875 Table 1 List of the climate and basin attributes used in this study.

876

877 Table 2 Median performance of the models trained under our proposed framework
878 (baseline models in the parenthesis) for two training scenarios.

879

880

881

882

883

884

885

886

887

888

889

890

891

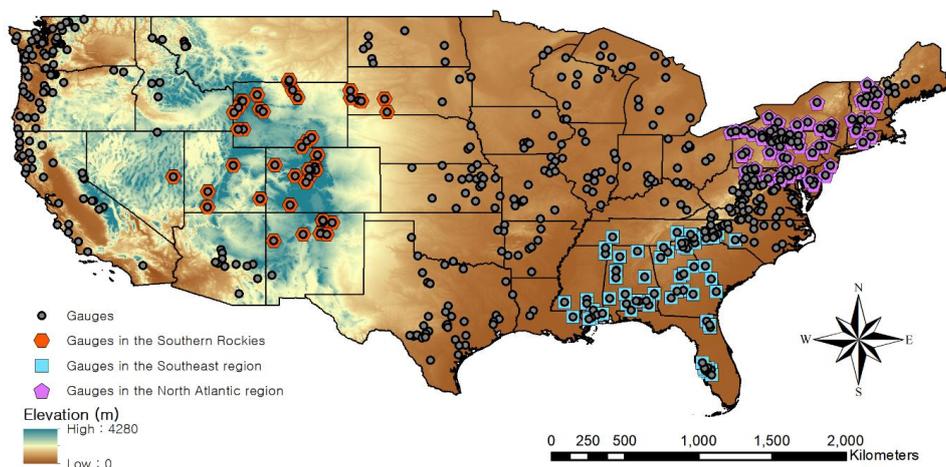
892



893

894

895



896

897 Figure 1 Locations of the selected 531 gauges in the United States. The gauge stations across
898 three homogenous regions are also illustrated.

899

900

901

902

903

904

905

906

907

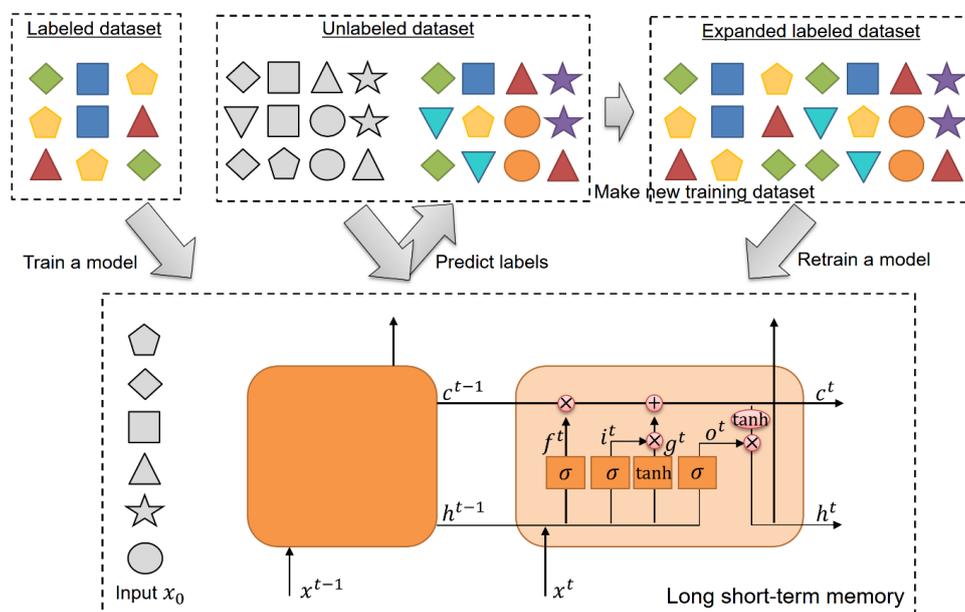
908



909

910

911



912

913 Figure 2 Overview of the semi-supervised learning-based framework proposed in this study.
 914 Here, the colored geometric shapes represent the recorded samples, pairing climate forcing data
 915 with streamflow.

916

917

918

919

920

921

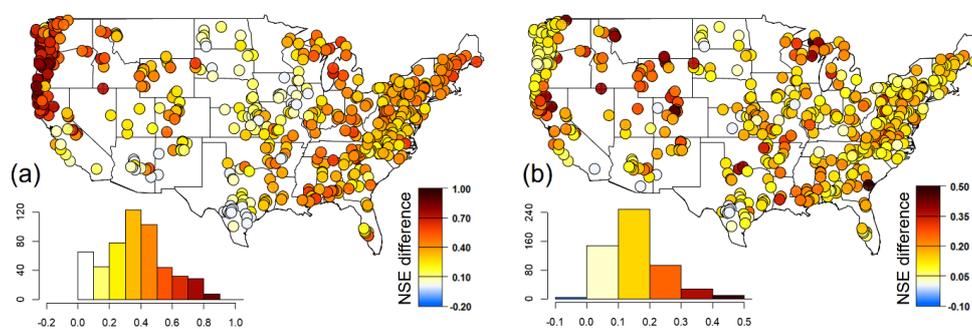
922



923

924

925



926

927 Figure 3 Difference of NSE results of idv-LSTM compared to their baseline models for (a)
928 single and (b) multi-year training scenarios. The color maps are limited to [-0.2, 1.0] for the
929 single training scenario and [-0.1, 0.5] for the multi-year training scenario for enhanced
930 visualization.

931

932

933

934

935

936

937

938

939

940

941

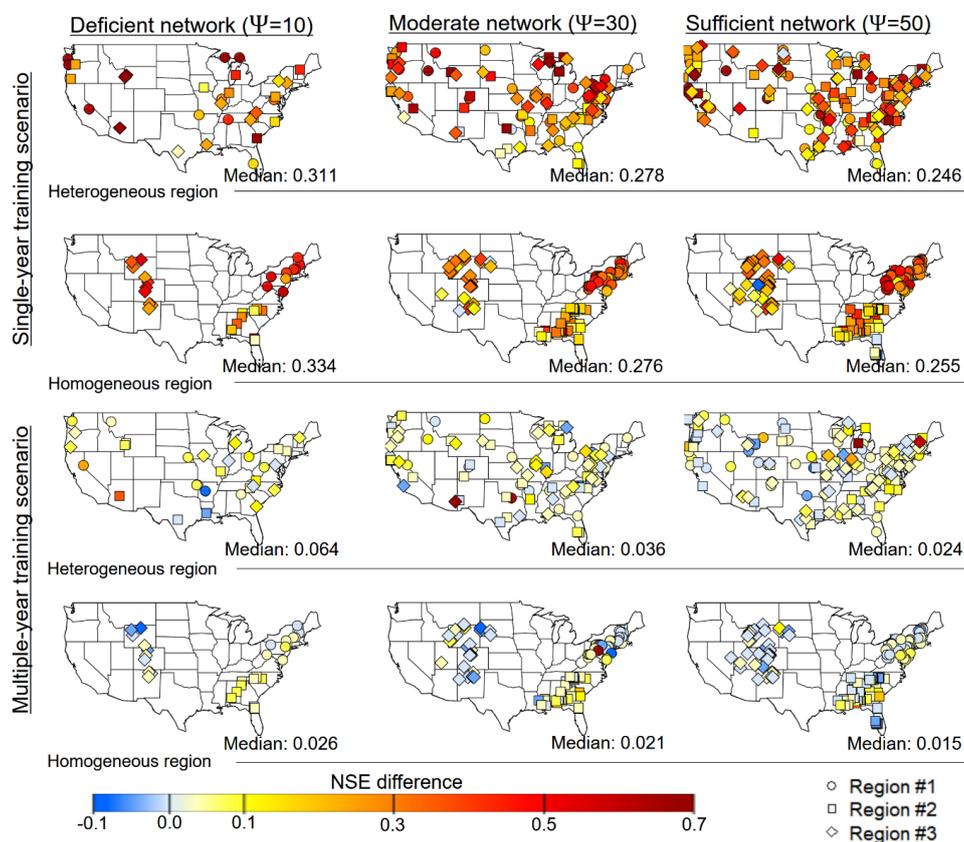
942



943

944

945



946

947 Figure 4 Difference of NSE results of rgn-LSTM compared to their baseline models across
 948 experimental factors including three defined regions, two training scenarios, and three basin
 949 densities in network. Here, the median NSE differences across basins in three defined regions
 950 are presented in each plot.

951

952

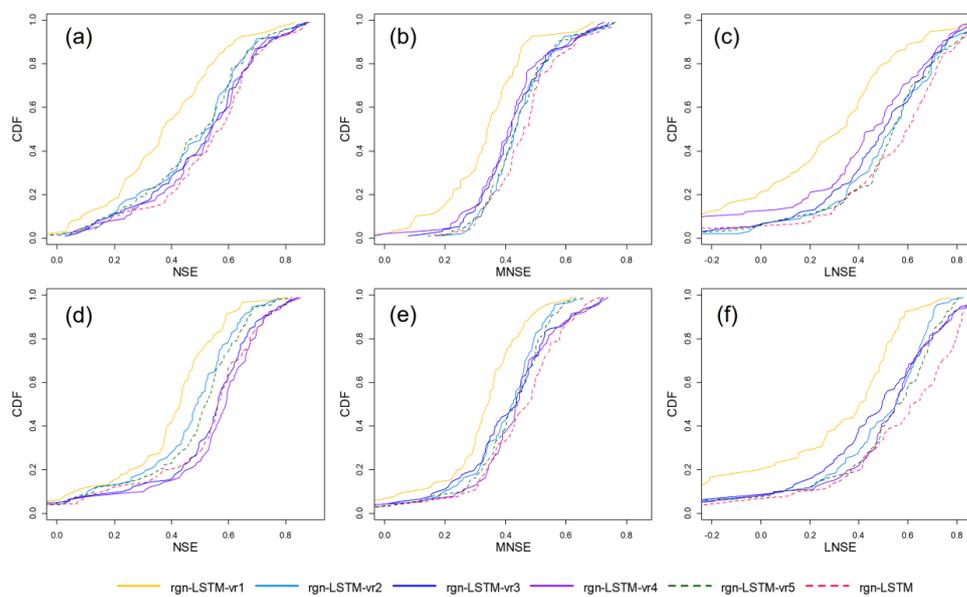
953



954

955

956



957

958 Figure 5 Cumulative density functions of the results of the annealing process on rgn-LSTM
959 with the single-year training scenarios obtained for basins across (a), (b), (c) heterogeneous
960 region; and (d), (e), (f) homogeneous region. Here, three metrics, namely NSE (first column),
961 MNSE (second column), and LNSE (last column), are utilized.

962

963

964

965

966

967

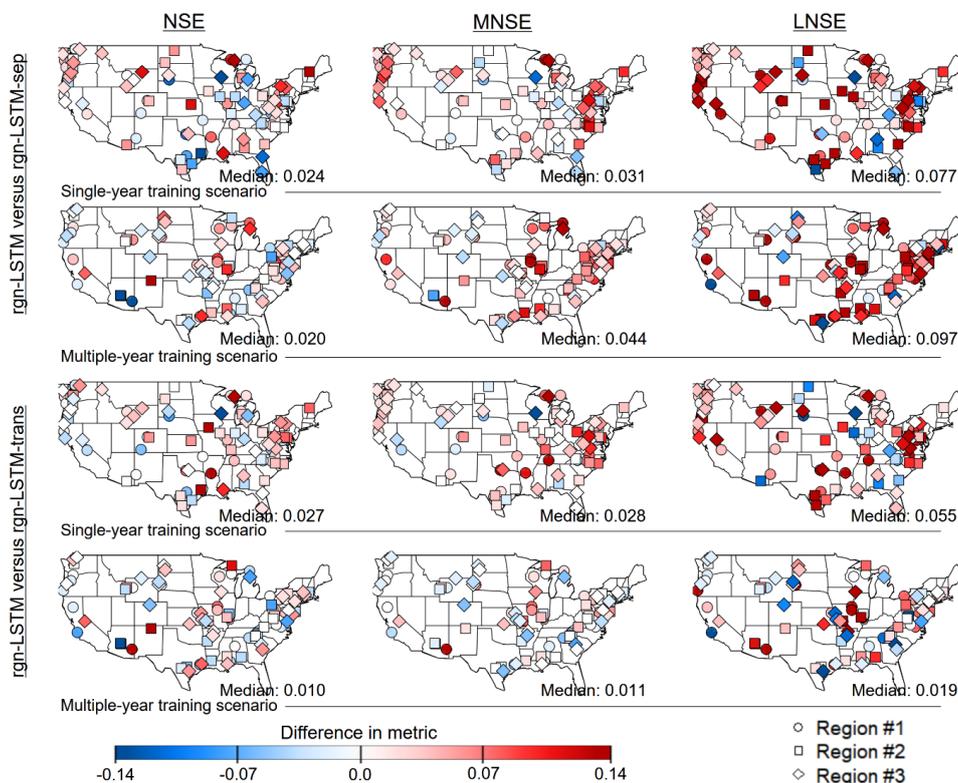
968



969

970

971



972

973 Figure 6 Difference of performance in the three metrics, NSE (first column), MNSE (second
 974 column), and LNSE (third column), of rgn-LSTM compared to the two fine-tuning approaches
 975 (rgn-LSTM-sep and rgn-LSTM-trans) across three basin networks in heterogeneous regions.
 976 Here, the median NSE differences across basins in three defined regions are presented in each
 977 plot.

978

979

980

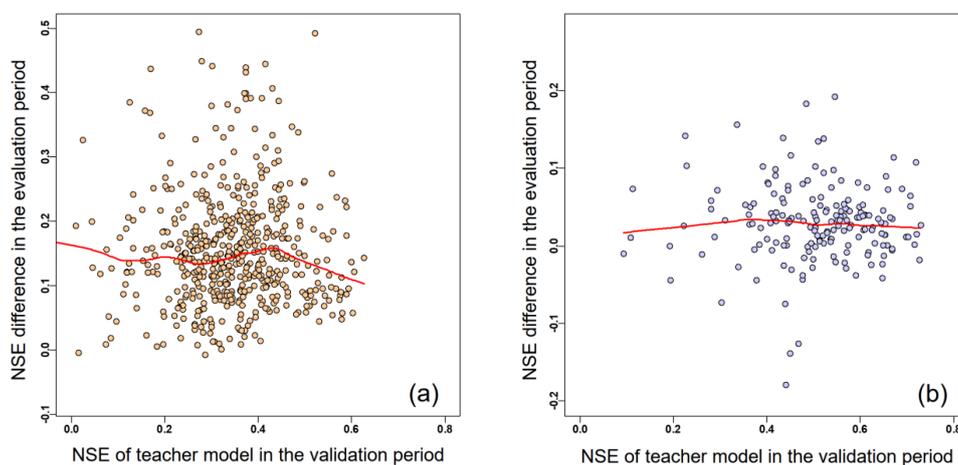
981



982

983

984



985

986 Figure 7 Plots comparing the NSEs achieved by the teacher model during the validation period
987 and the improved NSEs between the student models and baseline models during the evaluation
988 period under two model setting: (a) idv-LSTM and (b) rgn-LSTM with the moderate density
989 network, in multi-year training scenarios. Additionally, the expected improved NSEs
990 corresponding to the NSE of the teacher model are depicted, along with a Lowess fit
991 represented by a red line.

992

993

994

995

996

997

998

999



1000

1001

1002 Table 1 List of the climate and basin attributes used in this study.

| | Variable name | Description | Unit |
|---------------------|----------------------|---|------------------|
| | PRCP | Precipitation | mm |
| Climate forcing | T_{\max} | Maximum air temperature | °C |
| | T_{\min} | Minimum air temperature | °C |
| | SRAD | Short-wave radiation | W/m ² |
| | VP | Vapor pressure | Pa |
| | p_mean | Catchment mean daily precipitation | mm |
| | pet_mean | Catchment mean daily potential evapotranspiration | mm |
| | p_seasonality | Seasonality and timing of precipitation | - |
| | frac_snow | Fraction of precipitation falling as snow | - |
| | aridity | Ratio of catchment mean PET to mean precipitation | - |
| | high_prec_freq | Frequency of high precipitation days ($\geq 5 \times p_{\text{mean}}$) | Days |
| | high_prec_dur | Average duration of high precipitation events (number of consecutive days $\geq 5 \times p_{\text{mean}}$) | Days |
| | low_prec_freq | Frequency of dry days (< 1 mm/day) | Days |
| | low_prec_dur | Average duration of dry periods (number of consecutive days < 1 mm/day) | Days |
| Basin attributes | soil_depth_pelletier | Depth to bedrock (maximum 50m) | m |
| | soil_depth_statsgo | Soil depth (maximum 1.5m) | m |
| | soil_porosity | Volumetric porosity | - |
| | soil_conductivity | Saturated hydraulic conductivity | cm/hr |
| | max_water_content | Maximum water content of the soil | m |
| | sand_frac | Fraction of sand | % |
| | silt_frac | Fraction of silt | % |
| | clay_frac | Fraction of clay | % |
| | carbonate_rocks_frac | Fraction of the catchment area characterized as "Carbonate sedimentary rocks" | % |
| | geol_permeability | Subsurface permeability (log10) | - |
| | elev_mean | Catchment mean elevation | m |
| | slope_mean | Catchment mean slope | m/km |
| | area_gauges | Catchment area | km ² |
| | frac_forest | Forest fraction | % |
| | lai_max | Maximum monthly mean of the leaf area index | - |
| | lai_diff | Difference between the maximum and minimum monthly mean of | - |



1003
 1004

| | | | |
|----------|--|--|---|
| | | the leaf area index | |
| gvf_max | | Maximum monthly mean of the green vegetation fraction | - |
| gvf_diff | | Difference between the maximum and minimum monthly mean of the green vegetation fraction | - |

Table 2 Median performance of the models trained under our proposed framework (baseline models in the parenthesis) for two training scenarios.

| Mod el | No. of basins | Train ing scena rio | NSE _{median} | | MNSE _{median} | | LNSE _{median} | |
|------------------|--|---------------------------------|-----------------------------|---------------------------|-----------------------------|---------------------------|-----------------------------|---------------------------|
| | | | Heteroge neous region | Homogen eous region | Heteroge neous region | Homogen eous region | Heteroge neous region | Homogen eous region |
| idv- LST M | 531 basins | Single -year trainin g | 0.392 (0.002) | | 0.314 (0.076) | | 0.047 (-0.193) | |
| | | Multi- year trainin g | 0.576 (0.411) | | 0.441 (0.301) | | 0.418 (0.044) | |
| rgn- LST M | 30 basins in 3 experim ental trials | Single -year trainin g | 0.521 (0.175) | 0.531 (0.145) | 0.452 (0.206) | 0.421 (0.176) | 0.491 (- 0.634) | 0.624 (0.070) |
| | | Multi- year trainin g | 0.619 (0.563) | 0.650 (0.627) | 0.482 (0.446) | 0.536 (0.481) | 0.595 (0.461) | 0.700 (0.620) |
| rgn- LST M | 90 basins in 3 experim ental trials | Single -year trainin g | 0.570 (0.213) | 0.538 (0.241) | 0.483 (0.259) | 0.463 (0.226) | 0.580 (- 0.312) | 0.609 (0.022) |
| | | Multi- year trainin g | 0.667 (0.608) | 0.689 (0.673) | 0.549 (0.492) | 0.553 (0.515) | 0.666 (0.577) | 0.755 (0.672) |
| rgn- LST M | 150 basins in 3 | Single -year trainin g | 0.567 (0.263) | 0.542 (0.246) | 0.506 (0.241) | 0.467 (0.230) | 0.631 (0.166) | 0.636 (0.135) |



| | | | | | | | |
|---------|---------|---------|---------|---------|---------|---------|---------|
| experim | Multi- | | | | | | |
| ental | year | 0.675 | 0.704 | 0.555 | 0.579 | 0.691 | 0.750 |
| trials | trainin | (0.638) | (0.676) | (0.532) | (0.543) | (0.645) | (0.668) |
| | g | | | | | | |

1005