

# Exploring Variable Synergy in Multi-Task Deep Learning for Hydrological Modeling

Wenyu Ouyang<sup>1</sup>, Xuezhi Gu<sup>1</sup>, Lei Ye<sup>1</sup>, Xiaoning Liu<sup>1</sup>, and Chi ZHANG<sup>1</sup>

<sup>1</sup>Dalian University of Technology

November 22, 2023

## Abstract

Despite advances in hydrological Deep Learning (DL) models using Single Task Learning (STL), the intricate relationships among multiple hydrological components and model inputs might not be comprehensively encapsulated. This study employed a Long Short-Term Memory (LSTM) neural network and the CAMELS dataset to develop a Multi-Task Learning (MTL) model, predicting streamflow and evapotranspiration across multiple basins. An optimal multi-task loss weight ratio was determined manually during the validation phase for all 591 selected basins with streamflow data-gaps under 5%. During test period, MTL showed median Nash-Sutcliffe Efficiency predictions for streamflow and evapotranspiration at 0.69 and 0.92, consistent with two STL models. The MTL's strength appeared when predicting the non-target variable, surface soil moisture, using probes derived from LSTM cell states—representative of the internal DL model workings. This prediction showed a median correlation coefficient of 0.90, surpassing the 0.88 and 0.89 achieved by the streamflow and evapotranspiration STL models, respectively. This outcome suggests that MTL models could reveal additional rules aligned with hydrological processes through the inherent correlations among multiple hydrological variables, thereby enhancing their reliability. We termed this as “variable synergy,” where MTL can simultaneously predict varied targets with comparable STL performance, augmented by its robust internal representation. Harnessing this, MTL promises enhanced predictions for high-cost observational variables and a comprehensive hydrological model.

## Hosted file

977935\_0\_art\_file\_11538232\_s3bxwt.docx available at <https://authorea.com/users/695536/articles/684181-exploring-variable-synergy-in-multi-task-deep-learning-for-hydrological-modeling>

## Hosted file

977935\_0\_supp\_11538212\_s3bwnn.docx available at <https://authorea.com/users/695536/articles/684181-exploring-variable-synergy-in-multi-task-deep-learning-for-hydrological-modeling>

1 **Exploring Variable Synergy in Multi-Task Deep Learning for Hydrological Modeling**

2 **Wenyu Ouyang<sup>1</sup>, Xuezhi Gu<sup>1</sup>, Lei Ye<sup>1</sup>, Xiaoning Liu<sup>1</sup>, Chi Zhang<sup>1</sup>**

3 <sup>1</sup> School of Hydraulic Engineering, Dalian University of Technology, Dalian, China

4 Corresponding author: Lei Ye (yelei@dlut.edu.cn)

5 **Key Points:**

- 6
- Multi-task learning matched single-task models in spatiotemporal extrapolation accuracy
  - 7 • Reliability of multi-task learning is proven by enhanced correlation in probe predictions
  - 8 • Termed "variable synergy" for multi-task learning, highlighting its superior modeling.
- 9

## 10 **Abstract**

11 Despite advances in hydrological Deep Learning (DL) models using Single Task  
12 Learning (STL), the intricate relationships among multiple hydrological components and model  
13 inputs might not be comprehensively encapsulated. This study employed a Long Short-Term  
14 Memory (LSTM) neural network and the CAMELS dataset to develop a Multi-Task Learning  
15 (MTL) model, predicting streamflow and evapotranspiration across multiple basins. An optimal  
16 multi-task loss weight ratio was determined manually during the validation phase for all 591  
17 selected basins with streamflow data-gaps under 5%. During test period, MTL showed median  
18 Nash-Sutcliffe Efficiency predictions for streamflow and evapotranspiration at 0.69 and 0.92,  
19 consistent with two STL models. The MTL's strength appeared when predicting the non-target  
20 variable, surface soil moisture, using probes derived from LSTM cell states—representative of  
21 the internal DL model workings. This prediction showed a median correlation coefficient of  
22 0.90, surpassing the 0.88 and 0.89 achieved by the streamflow and evapotranspiration STL  
23 models, respectively. This outcome suggests that MTL models could reveal additional rules  
24 aligned with hydrological processes through the inherent correlations among multiple  
25 hydrological variables, thereby enhancing their reliability. We termed this as "variable synergy,"  
26 where MTL can simultaneously predict varied targets with comparable STL performance,  
27 augmented by its robust internal representation. Harnessing this, MTL promises enhanced  
28 predictions for high-cost observational variables and a comprehensive hydrological model.

## 29 **1 Introduction**

30 Deep learning (DL) models, specifically Long Short-Term Memory (LSTM) neural  
31 networks (Hochreiter & Schmidhuber, 1997), have exhibited notable proficiency for data  
32 integration and generalization in hydrological modeling (Feng et al., 2020; Kratzert, Klotz,  
33 Herrnegger, et al., 2019; Kratzert, Klotz, Shalev, et al., 2019; Ma et al., 2021). Their ability to  
34 efficiently leverage big data, discern high-dimensional relationships between variables and  
35 building general models, as posited by the Universal Approximation Theorem (Hornik et al.,  
36 1989), has been noteworthy in hydrology (Nearing et al., 2021; Shen, 2018). Consequently, they  
37 were widely employed in modeling and predicting a range of hydrological variables, such as  
38 streamflow, soil moisture, water temperature and dissolved oxygen (Liu et al., 2022; Nearing et  
39 al., 2021; Rahmani et al., 2021; Zhi et al., 2023). Despite these advancements, DL models might

40 learn improper patterns in hydrological modeling, even in the presence of robust goodness-of-fit  
41 results (Yokoo et al., 2022). A possible reason for this could be the focus of many deep-learning-  
42 based hydrological models on univariate modeling, which means they center their simulations on  
43 a single variable. Such an approach might increase the risk of overfitting in single-variable  
44 modeling, leading to an inadequate representation of relationships between model inputs and  
45 various hydrological components.

46 Conventionally, Physically Based Hydrological Models (PBHM) are also calibrated  
47 primarily using single variable data, commonly streamflow (Herman et al., 2018). However,  
48 some research has emphasized that models calibrated exclusively with streamflow may generate  
49 inadequate simulations for other water balance components (Becker et al., 2019; Tobin &  
50 Bennett, 2017; Yassin et al., 2017). Given that the hydrological process encompasses a multitude  
51 of variables involved in complex physical subprocesses, including surface and subsurface  
52 streamflow, soil water, and evapotranspiration (Shah et al., 2021), it is reasonable to incorporate  
53 additional components in the calibration of hydrological models. This could aid in constraining  
54 the solution of model parameters within a more viable parameter space (Dembélé, Hrachowitz, et  
55 al., 2020). Previous studies in physics-based modeling have shown that by incorporating a more  
56 rational representation of hydrological processes and calibrating model parameters with multiple  
57 model outputs, the overall predictive accuracy of hydrological variables could be improved, both  
58 in temporal and spatial generalization (Dembélé, Ceperley, et al., 2020; Tong et al., 2021, 2022).

59 While PBHMs are often calibrated using single-variable data, it is essential to note that  
60 they inherently consider the physical mechanisms of all involved variables through meaningful  
61 equations. Therefore, despite potential imperfection, PBHMs generally exhibit a reduced  
62 tendency for overfitting. On the other hand, due to their layered design and flexible architecture,  
63 DL models are more vulnerable to overfitting for one target. For example, many PBHMs can  
64 reasonably estimate evapotranspiration (Dembélé, Hrachowitz, et al., 2020; Shah et al., 2021;  
65 Yeste et al., 2023), whereas DL models struggle to predict it without direct training. This  
66 underscores the importance of further research into multi-variable calibration. Moreover, the  
67 advancements in hydrological remote sensing have facilitated the accumulation of extensive  
68 remotely-sensed hydrological variable data (McCabe et al., 2017), forming a basis for the  
69 exploration and analysis of DL models with multiple interrelated outputs.

70 In machine learning, multi-task learning (MTL) is an approach that enables a model to  
71 simultaneously learn the relationships between inputs and outputs of multiple tasks (Zhang &  
72 Yang, 2022). To learn the shared information between different tasks, the model needs to  
73 establish connections between the parameter spaces of different tasks in the MTL model. Hard  
74 Parameter Sharing is a prevalent method for achieving MTL (Vandenhende et al., 2022). This  
75 approach allows multiple tasks to share some encoding layers, known as shared layers, along  
76 with different task-specific layers for decoding and output. This method allows the MTL model  
77 to simultaneously learn correlations between multiple tasks and the unique features intrinsic to  
78 each task. Shared layers minimize memory usage during operation and eliminate computational  
79 cost of features within the shared layers, thereby improving the efficiency of training and testing  
80 relative to multiple single-task learning (STL) models (Vandenhende et al., 2019). Furthermore,  
81 the complementary information shared among related tasks may enable the model to learn a  
82 more generalized function relationship (Standley et al., 2020), thereby reducing the risk of  
83 overfitting.

84 Given the intrinsic interconnectedness of hydrological variables within a water cycle  
85 process, it is plausible to introduce MTL to hydrological deep learning-based modeling. Several  
86 studies have started to investigate the efficacy of MTL in hydrological models. Initial studies in  
87 MTL hydrological modeling primarily focused on incorporating more components in water  
88 balance, especially at large scales with abundant data. These studies utilized the water balance  
89 equation as a physical constraint and ensure that multiple interrelated hydrological processes are  
90 jointly optimized (Kraft et al., 2020). At the basin scale, Sadler et al. (2022) undertook research  
91 on MTL modeling in daily streamflow and water temperature, revealing that for certain sites and  
92 some MTL settings (like the scaling factor, denoting the ratio of loss from different tasks), MTL  
93 could enhance prediction accuracy across multiple tasks. Li et al. (2023) improved streamflow  
94 modeling with spatiotemporal DL models and an MTL approach in three basins. Building on  
95 these advancements, MTL has been adapted for a variety of hydrological targets, such as soil  
96 moisture (satellite and local in situ) (Liu et al., 2023), satellite precipitation estimation (rain/no-  
97 rain classification and rain rate) (Bannai et al., 2023), and aquifer transmissivity and storativity  
98 (Vu & Jardani, 2022). However, as the trend of modeling multiple variables has emerged, the  
99 precise benefits of MTL and how it behaves in terms of temporal and spatial generalization still  
100 not be fully understood, particularly in scenarios with large-sample basins.

101 Deep learning models have the potential to revolutionize our understanding of  
102 hydrological processes, but their reliability remains a topic of ongoing research. This study  
103 aimed not only to assess a model's potential for enhancing predictive performance but also to  
104 gauge its reliability by verifying if the input-output correlations learned by the model align with  
105 the established laws of hydrological processes. Various interpretative methods exist for  
106 hydrological deep learning models (Hu et al., 2021; Kratzert et al., 2021; Schmidt et al., 2020),  
107 but most are primarily used to analyze the attribution of input variables, not the internal states of  
108 DL models, thus posing challenges for our objectives. In natural language processing, supervised  
109 models known as “probes” have been devised to predict properties from representations  
110 (Belinkov et al., 2017; Hewitt & Liang, 2019), offering a way to inspect the learnt patterns of  
111 deep learning models. By leveraging such interpretability methods, we aimed to discern what the  
112 DL model truly learns for hydrological modeling. For example, Lees et al. (2022) adopted the  
113 probe method to analyze the relationship between the cell state and a non-target hydrological  
114 variable, thereby examining the plausibility of the processes LSTMs acquired.

115 The aim of this study was to construct MTL deep neural network models and conduct a  
116 comprehensive evaluation of their predictive performance in terms of both temporal and spatial  
117 generalizability across large-sample basins. We also investigated whether these models could  
118 learn more dependable correlations, potentially providing new “correct” insights that align with  
119 hydrological laws. Our approach integrated both MTL and STL techniques and evaluated their  
120 performance to examine the generalization capabilities and overall reliability of MTL.

## 121 **2 Data and Methods**

122 We utilized the Catchment Attributes and Meteorology for Large-Sample Studies  
123 (CAMELS) dataset (Addor et al., 2017), comprising 671 relatively undisturbed basins across the  
124 contiguous United States (CONUS). This dataset was widely used as a benchmark dataset for  
125 hydrological deep-learning-based modeling (Fang et al., 2022; Feng et al., 2020; Jiang et al.,  
126 2020; Kratzert, Klotz, Shalev, et al., 2019; X. Li et al., 2022). Given that evapotranspiration is  
127 frequently observed via remote sensing (Xu et al., 2019) and can serve as an output variable in  
128 hydrological models (Zhao, 1992), we integrated remotely observed evapotranspiration with the  
129 CAMELS dataset and considered it along with ground-observed streamflow as MTL target.  
130 Initially, we evaluated the performance of the models using the CAMELS dataset. As LSTM

131 hydrological models can store hidden information that represents hydrological knowledge, we  
132 employed physical interpretability methods (Lees et al., 2022) to compare the correlation  
133 relationships between internal states and outputs of MTL and STL models. This approach  
134 facilitated an understanding of the internal states of the trained models and provided evidence for  
135 varying degrees of reliability among the models.

## 136 2.1 Dataset

137 The development of MTL models that simultaneously consider the output of multiple  
138 hydrological variables necessitates the assembly of datasets incorporating several hydrological  
139 model output variables. In the CAMELS dataset, for general hydrological modeling outputs, only  
140 streamflow data are obtained from observations, with the other variables' outputs derived through  
141 hydrological simulation. Thus, to explore the potential for MTL based on the CAMELS dataset,  
142 we expanded the available data for the basins in CAMELS.

143 We retrieved the evapotranspiration data from the MOD16A2 data product (Running et  
144 al., 2017) from the Google Earth Engine (GEE) data catalog (Gorelick et al., 2017). This dataset,  
145 comprising an 8-day temporal resolution and a 500-meter spatial resolution, spans from 2001-01-  
146 01 to the present. However, it's noteworthy that while most data collection periods are 8 days, the  
147 final collection period of each year is adjusted to 5 days for non-leap years and 6 days for leap  
148 years. The algorithm behind the MOD16 data product employs the Penman-Monteith equation,  
149 which is supplemented with daily meteorological reanalysis data and other MODIS remote  
150 sensing data products. The output includes several raster data layers, including actual  
151 evapotranspiration (ET). The pixel values in the ET data layer denote the sum of daily values for  
152 each resolution period. In this study, ET was used as the output observation for model training.  
153 To derive the basin-mean daily time series of ET data, we used Map-Reduce functions in GEE  
154 (Gorelick et al., 2017). Specifically, each pixel from the gridded ET data was allocated to a  
155 specific region, leveraging weighted reducers to ensure accurate assignment.

156 The CAMELS data covers the period from 1980-01-01 to 2014-12-31, while the ET data  
157 is only available from 2001-01-01 onwards. To secure a sufficient data period, we extended each  
158 time series in the supplementary CAMELS dataset to 2021-09-30. As a result, the period  
159 considered for all models in this study was from 2001-01-01 to 2021-09-30. The NLDAS-2  
160 (NASA, 2018), Phase 2 of the North American Land Data Assimilation System and one of the

161 sources of meteorological data in CAMELS, was used as the basin meteorological forcing data.  
 162 The daily time series basin-mean forcing data were obtained in GEE using the same method as  
 163 for ET. Additionally, we supplemented the streamflow data from 2015-01-01 to 2021-09-30  
 164 from the U.S. Geological Survey (USGS) National Water Information System (NWIS) (USGS,  
 165 2019). To mitigate the influence of excessive missing data on the results, we selected basins with  
 166 a streamflow data loss rate of less than 5% during the overall period analyzed. This resulted in  
 167 the inclusion of 591 of the 671 CAMELS basins. Attributes related to soil, geology, topography,  
 168 land use types, and climate from the CAMELS dataset were also used as inputs for all models in  
 169 this study. For more information on these inputs, see Table 1.

170 This study used surface soil moisture (SSM) to assess the reliability of the STL and MTL  
 171 models (details provided in section 2.5). The data source was the SMAP global SSM dataset  
 172 (Mladenova et al., 2019) from GEE. The basin-averaged SMAP grid data was obtained using the  
 173 same method as for ET data acquisition. Consequently, the daily time series data for the SSM of  
 174 each basin were compiled.

175 **Table 1.** Hydrological variables selected as inputs and outputs to single- and multi-task deep  
 176 learning models based on the augmented CAMELS dataset.

Variable Type	Variable Name	Description	Unit	
Forcings	total_precipitation	Daily total precipitation	kg/m <sup>2</sup>	
	potential_evaporation	Potential evaporation	kg/m <sup>2</sup>	
	temperature	Air temperature at 2 meters above the surface	°C	
	specific_humidity	Specific humidity at 2 meters above the surface	kg/kg	
	shortwave_radiation	Surface downward shortwave radiation	W/m <sup>2</sup>	
	potential_energy	Convective available potential energy	J/kg	
Attributes	elev_mean	Basin mean elevation	m	
	Terrain	slope_mean	Basin mean slope	m/km
		area_gages2	Basin area	km <sup>2</sup>
Land Cover	frac_forest	Forest proportion	-	
		lai_max	Maximum monthly mean of leaf area index	-
		lai_diff	Difference between the maximum and	-

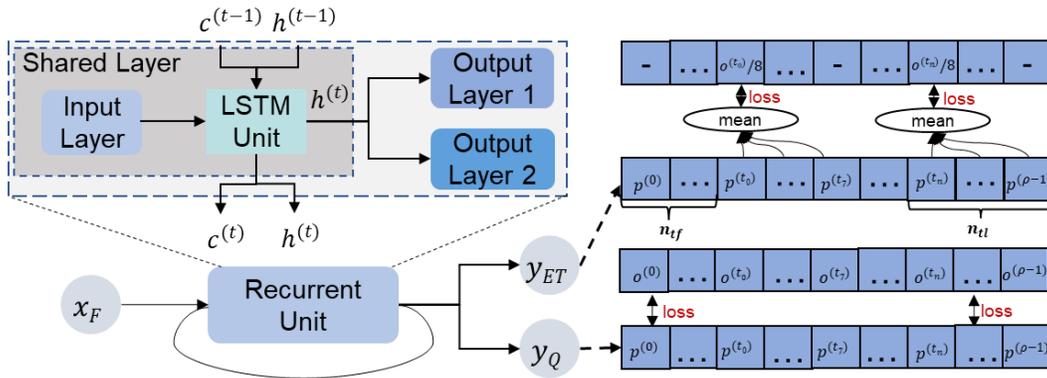
		minimum monthly mean values of the leaf area index	
	dom_land_cover_frac	Proportion of major land cover types to watershed area	-
	dom_land_cover	Major land cover types	-
Soil	root_depth_50	Average soil layer thickness containing the top 50% of the root system	m
	soil_depth_statgso	Soil depth	m
	soil_porosity	Soil porosity	-
	soil_conductivity	Saturated hydraulic conductivity	cm/hr
	max_water_content	Maximum soil water holding capacity	m
	Geology	geol_class_1st	Most common geological category in the watershed
geol_class_2nd		Second most common geological category in the watershed	-
geol_porosity		Subsurface porosity	-
geol_permeability		Subsurface permeability	m <sup>2</sup>
Model Outputs	streamflow	Daily streamflow in the outlet of a basin	m <sup>3</sup> /s
	evapotranspiration	Basin mean daily actual evapotranspiration	mm/day

## 177 2.2 Multi-task LSTM

178 LSTMs have become a prevalent choice in hydrological modeling due to their ability to  
179 capture temporal sequences and intricate patterns in the data. In this framework, MTL  
180 simultaneously optimizes multiple related tasks, leveraging shared representations, while STL  
181 focuses on optimizing one specific task. In this study, both MTL and STL models incorporated  
182 LSTM structures. The LSTMs, analogous to those proposed in prior research (Feng et al., 2020;  
183 Ma et al., 2021), including our previous study (Ouyang et al., 2021), operated as N-to-N models.  
184 This N-to-N term indicates that for every N input sequences, N output sequences are generated.  
185 These models leverage meteorological forcings and static basin attributes to predict daily  
186 discharge in the CAMELS dataset. We developed MTL models using a hard parameter sharing  
187 architecture (Vandenhende et al., 2022). The STLs' structure was totally same with the MTL's

188 except that they only calculated the loss of one output variable. With this setting, we controlled  
 189 the varying factor and the difference between STL and MTL was only the output.

190 Figure 1 presents the structure of a single time-step unit in the MTL hydrological model.  
 191 Central to this design is a shared layer, composed of a fully connected input layer and an LSTM  
 192 unit. The input layer consisted of two layers, each containing 256 neurons. As data progresses  
 193 through the LSTM, it attempts to express the intricate temporal patterns of hydrological  
 194 processes. Emerging from this shared space are multiple, parallel fully connected output layers,  
 195 each corresponding to a hydrological task. Each task-specific output's neurons were arranged in  
 196 two layers, with 128 and 1 neurons, respectively. Both the input and output fully connected  
 197 layers introduce non-linearity through the ReLU activation function. To summarize, during  
 198 forward computation of the model, inputs pass through a shared layer, generating long sequential  
 199 multiple feature variables. Then these variables are moved through different output layers to  
 200 generate the corresponding multiple outputs.



201  
 202 **Figure 1.** Illustration of the MTL hydrological model. The model inputs,  $\mathbf{x}_F$ , comprise a vector  
 203 of raw meteorological forcing inputs, and outputs,  $\mathbf{y}_Q$  and  $\mathbf{y}_{ET}$  represent the streamflow and  
 204 evapotranspiration, respectively. The LSTM's internal state in the  $t$ -th period is denoted by the  
 205 cell state,  $\mathbf{c}_t$  and hidden state,  $\mathbf{h}_t$ . In each period,  $\mathbf{p}_t$  represents the prediction and  $\mathbf{o}_t$  is the  
 206 observed data. The missing data in a given period is indicated by "-". The symbol "mean"  
 207 enclosed in a circle represents the mean value of selected periods.

208 Backpropagation in these models allows the independent updating of weight and bias  
 209 parameters for task-specific output layers, based solely on the losses of the current layer and  
 210 independent of other output losses. However, updates to the shared LSTM layer parameters  
 211 depend on multiple outputs. The following equations illustrate these updates:

$$212 \quad \theta_T(i+1) = \theta_T(i) - \alpha \nabla_{\theta_T} L(\theta_T(i), \theta_S(i))$$

$$213 \quad \theta_S(i+1) = \theta_S(i) - \alpha \sum_{j=1}^n \omega_j \left[ \nabla_{\theta_S}^{(j)} L(\theta_S(i), \theta_T^{(j)}(i)) \right]$$

214 In these equations,  $\theta$  denotes the weights and biases of the neural networks, with T and S  
 215 representing the task-specific output and shared layers, respectively. The index  $I$  signifies the  $i$ -th  
 216 training step,  $\alpha$  denotes the learning rate,  $\nabla$  represents the gradient of the loss function relative to  
 217 the weight parameter, and  $L(\cdot)$  is the loss function itself. The  $j$ -th specific task is represented by  $j$ ,  
 218  $\omega_j$  signifies the weights and bias corresponding to the  $j$ -th specific task and  $n$  stand for the total  
 219 number of tasks.

220 One of the challenges of constructing a multi-task learning model is balancing the loss  
 221 from each task. This balance is crucial to avoid one task from dominating the model training and  
 222 negatively impacting the learning of other tasks (Vandenhende et al., 2022). The MTL loss  
 223 function, represented by equation (3), calculates the overall loss value for all tasks, where  $L_{\text{MTL}}$   
 224 signifies the overall loss value for all tasks,  $L_j$  represents the loss value of the  $j$ -th task, and other  
 225 variables have the same meaning as in equation (2).

$$226 \quad L_{\text{MTL}} = \sum_{j=1}^n \omega_j \cdot L_j$$

$$227 \quad \sum_{j=1}^n \omega_j = 1$$

228 Balancing tasks can be achieved by setting task-specific weights, represented as  $\omega_j$ , in the  
 229 loss function. However, quantifying the weight of each task is challenging. Two usual  
 230 approaches to task balancing exist: the uncertainty weighting method (Cipolla et al., 2018) and  
 231 dynamic task prioritization (Guo et al., 2018). However, these methods adopt totally different  
 232 views on the significance of tasks. The former balances task losses by considering homoscedastic  
 233 uncertainty, assigning lesser weight to outputs with higher uncertainty and consequently higher  
 234 weight to simpler tasks. But the latter prioritizes the learning of difficult tasks by assigning them  
 235 higher task-specific weights.

236 A more direct and simpler approach is manual loss weight assignment, which was also  
 237 used in some related studies (B. Li et al., 2023; Sadler et al., 2022). This paper defined the loss  
 238 weight ratio  $\lambda$  as the ratio of evapotranspiration and streamflow variable loss weights,  $\frac{\omega_{ET}}{\omega_Q}$ .  
 239 During the training period, multiple  $\lambda$  values were assigned, each corresponding to an MTL

240 model trained for all basins simultaneously. The model demonstrating the best overall prediction  
241 performance during the validation period was chosen for testing.

242 The MTL model was designed to produce daily predictions for both streamflow and  
243 evapotranspiration. Although daily streamflow observation data is available, evapotranspiration  
244 observation data is cumulative and represents values over an 8-day interval. This interval is  
245 adjusted to 5 or 6 days in regular and leap years, respectively, to account for the final period of  
246 each year. Therefore, a specific design for the loss function calculation is necessary. As shown in  
247 Figure 1, the observed streamflow values were directly compared with predicted values.  
248 Meanwhile, predicted ET values were averaged over a period before being used to calculate the  
249 loss function. The first  $n_{tf}$  or last  $n_{tl}$  time-steps of the whole period could begin or end with a  
250 duration of less than 8 days. In such situations, we ignored the first  $n_{tf}$  non-value time-steps and  
251 multiplied the final observed value by  $n_{tl}/8$ ,  $n_{tl}/5$ , or  $n_{tl}/6$ , depending on whether the last period  
252 was the final period in a regular or leap year. Throughout the model training phase, the root-  
253 mean-square error (RMSE) acted as the loss function. This same RMSE metric was applied to  
254 calculate the loss functions for each individual output under the MTL mode.

### 255 2.3 General settings

256 All models employed in this paper utilized the same input variables, including 6  
257 meteorological forcing variables and 17 attribute variables pertinent to soil, geology, topography,  
258 land use type, and climate. The details are provided in Table 1. One of the distinct advantages of  
259 deep learning models is their ability to automatically extract input features from an end-to-end  
260 perspective, rather than manually analyzing and extracting features from multiple input  
261 variables. Hence, the basin attribute data were directly copied to each period and concatenated  
262 with the meteorological input, creating the model's input vector without necessitating manual  
263 selection.

264 The settings for data preprocess and model training aligned with our previous and related  
265 research (Ouyang et al., 2021; Rahmani et al., 2021) and were consistently applied to all models,  
266 including STL and MTL models, to ensure comparability.

267 Before model training, normalization of input and output data samples is essential for the  
268 efficient optimization of the neural network weight by the gradient descent algorithm during

269 subsequent training. Test data also require normalization, and the statistical data used for  
270 normalization during testing is that used for training. After the model completed its predictions,  
271 the results are re-normalized back to their original dimension.

272 Consistent with our previous research, the Adadelta algorithm, an adaptive learning rate  
273 scheme (Zeiler, 2012), was chosen as the optimization method for performing stochastic gradient  
274 descent on the neural network model parameters. To mitigate overfitting, dropout regularization  
275 was implemented during the training of LSTM models. Dropout applies a fixed mask, meaning  
276 once a connection weight is set to zero, it stays at zero for the entire training process. The loss  
277 function was the root mean square error between the observed and predicted values. The  
278 hyperparameter settings of all models in this study were as follows: the mini-batch size was 100,  
279 the training sequence length was 365, the number of hidden units per layer was 256, and the  
280 LSTM dropout rate was 0.5.

281 In the evaluation phase, the Nash-Sutcliffe Efficiency (NSE) score was employed to  
282 assess streamflow and evapotranspiration prediction. NSE is a metric particularly suited to  
283 evaluate hydrological predictions. Additionally, other common metrics, such as the mean  
284 difference between modeled and observed values (Bias), RMSE, and Pearson's correlation  
285 (Corr), were also used to evaluate the models.

## 286 2.4 Experiments

287 This study devised two experiments to ascertain the conditions under which an MTL  
288 model could enhance the simultaneous prediction of each variable compared to STL models. In  
289 experiment A, we partitioned the dataset into training, validation, and test sets, with the  
290 validation data assisting in finding the optimal multi-task loss weight ratio  $\lambda$ . The evaluation  
291 metrics of the STL and MTL models for each output were compared in this experiment. In  
292 experiment B, we further investigated the temporal and spatial generalization capabilities of the  
293 STL and MTL models using scaling curves to gain a deeper understanding of their differences.

### 294 Experiment A: Comparison of STL and MTL models utilizing the entire dataset

295 We first constructed an MTL model that predicted both streamflow and  
296 evapotranspiration. To assess any potential improvement in predictive capability, the  
297 performance of this model was compared with that of two STL models; one predicting

298 streamflow and the other predicting evapotranspiration. Notably, the STL model for streamflow  
299 did not encompass any input or output associated with evapotranspiration data, and vice versa.

300         Employing the multi-task balance strategy outlined in Section 2.2, the multi-task loss  
301 weight ratio  $\lambda$  was manually assigned. We chose five  $\lambda$  values (2, 1, 1/3, 1/8, and 1/24) to  
302 conduct prediction experiments with the MTL model. The optimal model for the testing period  
303 was identified by evaluating the NSE values achieved for each variable in the basins during the  
304 validation period. The training, validation, and test periods were from 2001-10-01 to 2011-09-30,  
305 2011-10-01 to 2016-09-30, and 2016-10-01 to 2021-09-30, respectively.

306         Experiment B: Assessment of model temporal and spatial generalization

307         Generally, supplying more data for DL models often leads to superior model  
308 performance. As the number of basins expands, the temporal and spatial generalization of the  
309 models usually improve. Scaling curves, which depict the behavior of scaling relative to the  
310 amount of training data (Tsai et al., 2021), could be used to analyze how the models behave as  
311 the number of trained basins increases. By comparing the STL and MTL models, the conditions  
312 under which MTL models outperform STL models could be identified.

313         For all models, a percentage of basins were randomly chosen for training, with the  
314 remaining basins used for temporal and spatial generalization evaluation. We chose 11  
315 percentage values: 5, 10, 20, 25, 33, 50, 66, 75, 80, 90 and 95. To mitigate geospatial bias, we  
316 ensured that each case included basins from every LEVEL-II ecoregion (Omernik & Griffith,  
317 2014), rendering them representative of the entire group. When the number of basins was  
318 limited, the selection process could introduce bias. Hence, we employed cross-validation to  
319 randomly select basins from the entire dataset repeatedly and computed the average median  
320 metric value across all cases as the result. The training period was the same as in Experiment A  
321 (2001-10-01 to 2011-09-30). No specific validation period was assigned as it was determined  
322 based on the best multi-task loss weight ratio  $\lambda$  obtained from Experiment A. The test period was  
323 from 2011-10-01 to 2021-09-30.

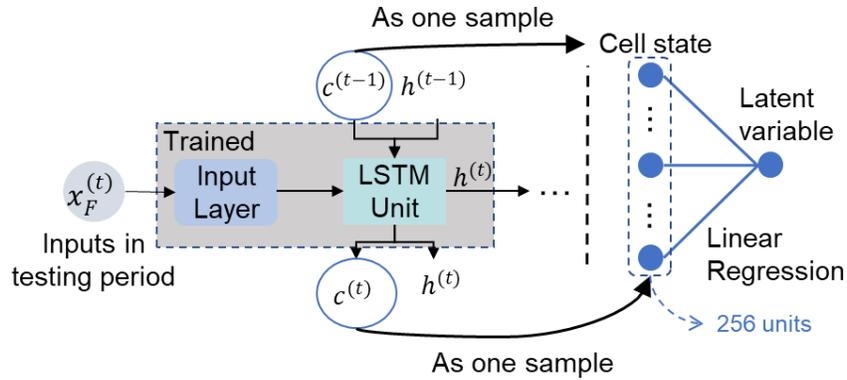
324         2.5 Reliability Assessment

325         We evaluated the reliability of deep learning models by comparing the predictive  
326 capabilities of probes in STL and MTL models. Our hypothesis posited that if the probes in MTL

327 models outperformed those in STL models in predicting non-target hydrological variables, then  
328 the MTL models were extracting more information, as the probes denoted the LSTM state  
329 vector's capacity to predict non-target variables. Such an outcome would further imply that MTL  
330 models were effectively identifying more credible correlations between inputs and multiple  
331 outputs. As DL models compress input information in their high-dimensional space based on the  
332 loss between observations and outputs, having more outputs implied that more information was  
333 encoded. Therefore, it was plausible to expect differences in the predictive performance of  
334 probes between STL and MTL models.

335         The implementation of the latent variable's probe is outlined in detail below. We began  
336 by training STL or MTL models, then we input the concatenated meteorological forcing data and  
337 attributes (XF) from the testing period into the trained models (as depicted in Figure 2). Next, we  
338 extracted cell states for all periods and use these to train a linear regression model. This model  
339 took 256 units from each sample over each period as input and generated a non-target variable as  
340 output. Subsequently, we produced predictions from the probe and compare them with  
341 observations. Notably, both the training and testing periods for the probe were included in the  
342 testing for both STL and MTL models. For this paper, we adopted a 4:1 ratio for the training-to-  
343 testing periods for the probes.

344         In both STL and MTL models used for streamflow (Q), evapotranspiration (ET), or both,  
345 surface soil moisture (SSM) was the non-target variable for STL or MTL models and serves as  
346 the target variable for the probe. In the STL model of streamflow, ET was the non-target  
347 variable, and for the STL model of ET, Q was the non-target variable. Even though the probe  
348 was typically used for non-target variables in deep learning models, we used it to probe both Q  
349 and ET in all STL and MTL models to thoroughly examine the differences between STL and  
350 MTL models. We evaluated the correlation of the predictive performance of probes in both STL  
351 and MTL models. The probes for ET or Q could assist us in understanding how probes behave  
352 with target variables. A stronger correlation for the SSM probe in the MTL model could suggest  
353 that the implicit information in the MTL model aligns more closely with the actual hydrological  
354 process.



355

356 **Figure 2.** Illustration of the training process of a latent variable's probe. Each cell state in one  
 357 period is considered a single sample for the linear regressor. The input for the regressor matches  
 358 the size of the cell state (256 units), while the output size is 1, representing a single latent  
 359 variable.

### 360 3 Results

#### 361 3.1 Prediction performance of MTL and STL models

362 Upon examining the performance of the two variables during the validation period,  $\lambda=1/3$   
 363 was chosen for the MTL model evaluation in testing period. Further details can be found in  
 364 Supporting Information Figure S1. In this section, we focus on the results for testing period.

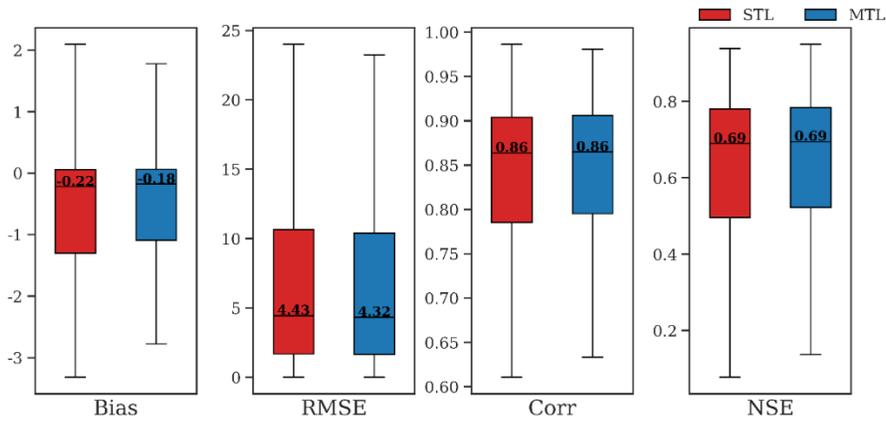
365 As depicted in Figure 3, the performance metrics of the MTL and STL models for both  
 366 streamflow and evapotranspiration prediction are relatively similar. For streamflow prediction,  
 367 the median value of the RMSE of the MTL model is 4.32 (m<sup>3</sup>/s), marginally lower than that of  
 368 the STL model. The Correlation and NSE median values are almost identical to those of the STL  
 369 model, at 0.86 and 0.69, respectively, albeit with slightly superior upper and lower boundaries of  
 370 the box plot. The Bias of the MTL model for streamflow prediction is nearer to 0 than that of the  
 371 STL model, showing an improvement of approximately 18%. The results for evapotranspiration  
 372 prediction, as displayed in Figure 3(b), follow a similar pattern with the RMSE, Correlation, and  
 373 NSE of the MTL model being -0.04 (mm/day), 0.96, and 0.92, respectively. These values are  
 374 equivalent to those of the STL model.

375 Previous studies calibrating physics-based hydrological models with multiple  
 376 hydrological output variables (Dembélé, Ceperley, et al., 2020; Dembélé, Hrachowitz, et al.,  
 377 2020; Tong et al., 2021) found that while using multiple outputs enhanced the simulation

378 accuracy for variables other than streamflow, the accuracy of the streamflow simulation itself  
 379 decreased. In contrast, our research highlighted that, when looking at the collective performance  
 380 across all basins, a deep-learning-based MTL model that not only slightly improved the  
 381 prediction performance of streamflow but also maintained the accuracy for evapotranspiration  
 382 prediction. Significantly, MTL models could simply consider multiple process and  
 383 simultaneously output multiple variables. Hence, in scenarios where various processes should be  
 384 considered, it was more reasonable to construct an MTL model rather than using multiple STL  
 385 models that only modeled one hydrological output variable at a time.

386

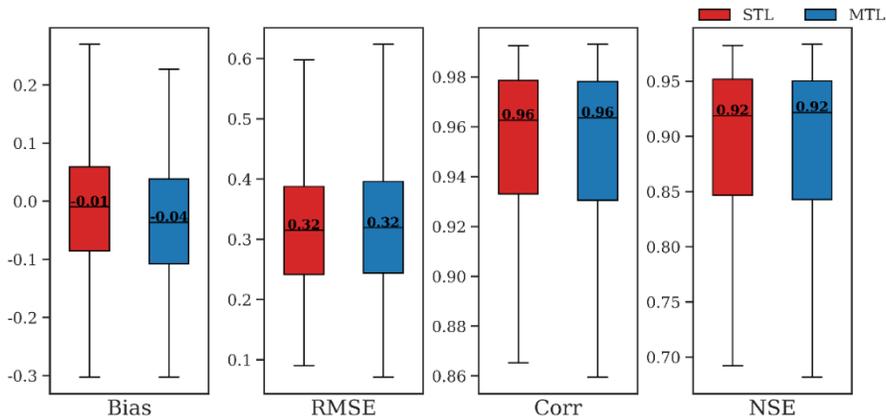
(a)



387

388

(b)



389

390 **Figure 3.** Statistical indicators of the streamflow and evapotranspiration prediction results of  
391 each STL model during the testing period and the MTL model under the  $\lambda=1/3$  scheme, which  
392 include Bias, RMSE, Corr, and NSE.

393 Figure 4 contrasts the performance of STL and MTL models for streamflow and  
394 evapotranspiration prediction across various basins. The NSE values for the STL and MTL  
395 models vary significantly among different basins. As illustrated in Figure 4(a), the integration of  
396 evapotranspiration into the MTL model doesn't invariably enhance streamflow prediction for all  
397 basins. About half of the basins (280) display superior streamflow predictions with the STL  
398 model, while the other 311 basins demonstrate improved predictions with the MTL model. Most  
399 basins are situated near the 1:1 line of equality, suggesting that the added dimension of  
400 evapotranspiration often led to subtle variations in prediction results in most basins. Figure 4(b)  
401 reveals that evapotranspiration prediction shows analogous patterns with minor differences  
402 between the STL and MTL models across most basins. Some basins demonstrate superior  
403 predictions with the MTL model, while others with the STL model.

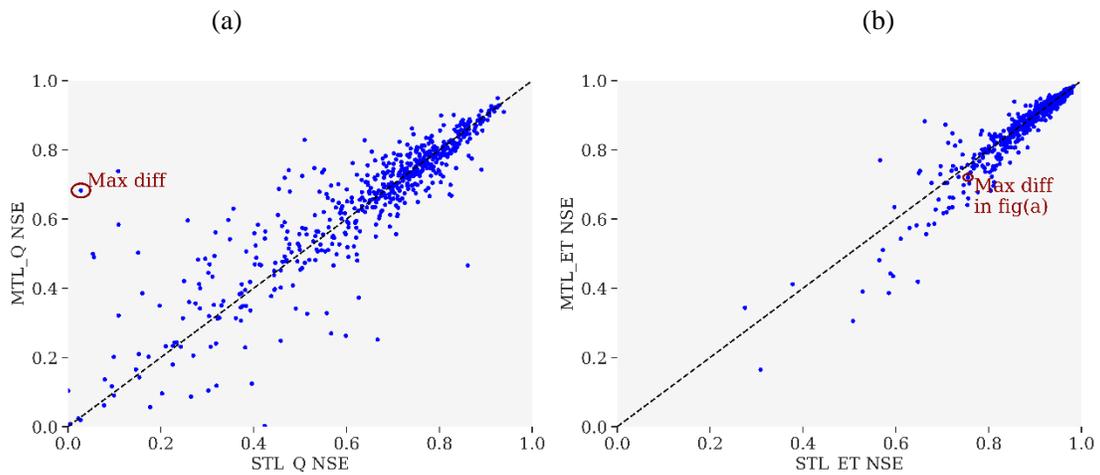
404 In some basins, particularly displayed in the top left corner of Figure 4(a), shown in the  
405 circle with label "Max diff", the differences in NSE for streamflow prediction are strikingly  
406 large. The MTL model exhibits a considerably higher NSE value for streamflow prediction of  
407 0.68 compared to the STL model's 0.03. However, the STL model performs more effectively in  
408 forecasting evapotranspiration in this basin, shown in Figure 4(b) with a circle label, with NSE  
409 values for STL and MTL models being 0.76 and 0.72, respectively.

410 Figure 4(c) and (d) feature a comparison between the streamflow and evapotranspiration  
411 prediction results of STL and MTL models, along with the observational data, in the basin where  
412 streamflow prediction saw the most significant improvement. The MTL model's streamflow  
413 prediction accurately sidesteps unrealistic high flow rates that don't align with observations,  
414 thereby enhancing the predictive performance. In terms of evapotranspiration prediction values,  
415 the MTL model is slightly lower than the STL model and exhibits limited consistency with the  
416 observation values.

417 These results suggested a competitive relationship among various output variables in  
418 MTL, where improving the prediction performance of one variable might lead to a decline in the  
419 performance for another. This phenomenon tied back to the concept of attaining a Pareto frontier  
420 in multi-objective optimization. Hydrological variables should ideally adhere to a single physical

421 law, regardless of apparent competitive relationships between multiple objectives, such as flood  
 422 control and power generation in reservoir operations. This competitive relationship indicated the  
 423 presence of latent variables influencing the formation processes of evapotranspiration and  
 424 streamflow that the current input failed to capture. Therefore, reaching the Pareto frontier of the  
 425 MTL process became crucial (Sener & Koltun, 2018).

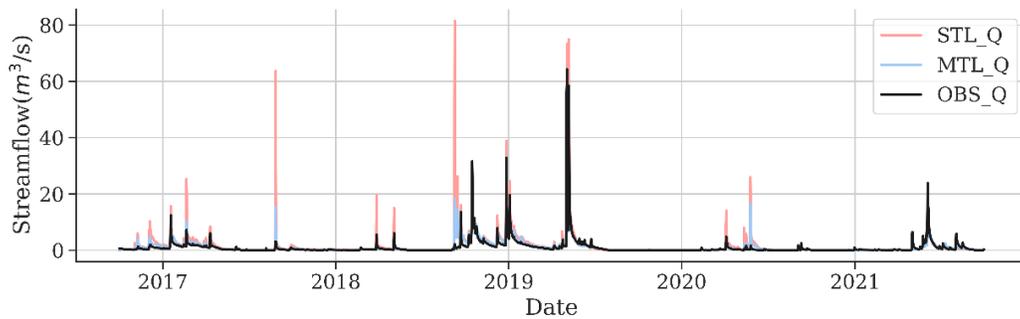
426



427

428

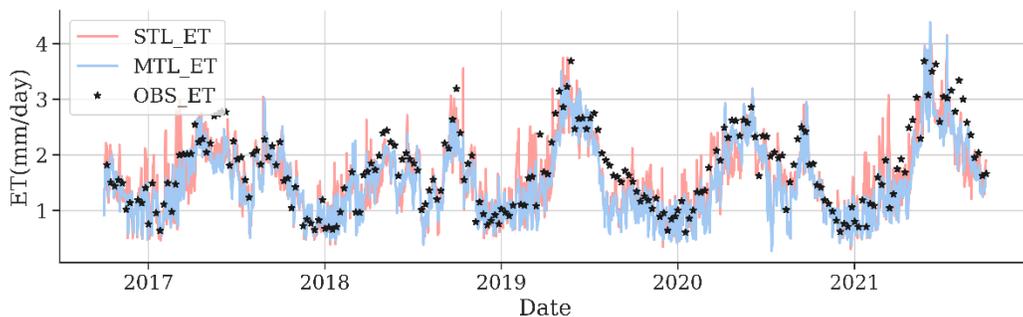
(c)



429

430

(d)



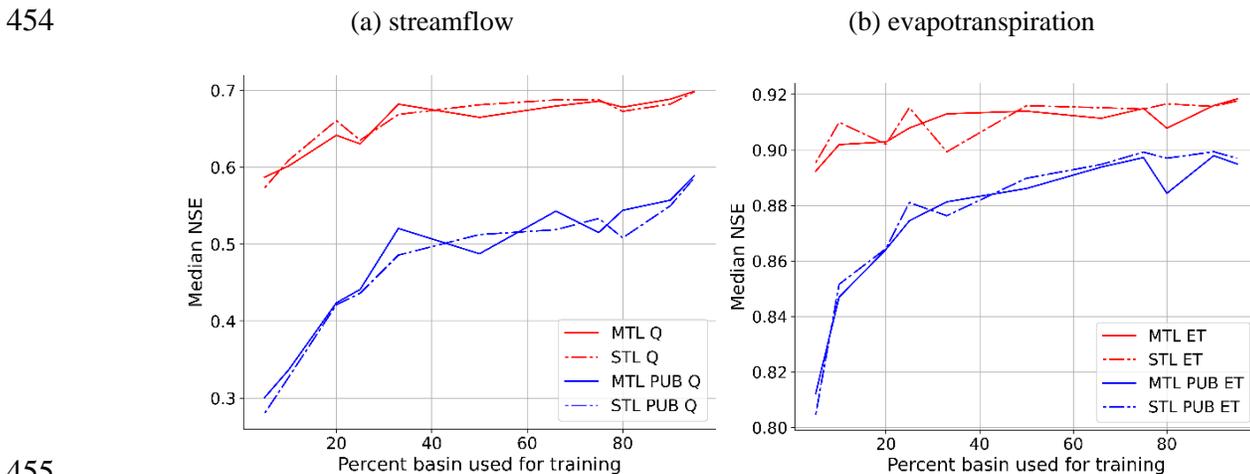
431

432 **Figure 4.** Streamflow and evapotranspiration predictions of STL and MTL models. During the  
 433 test period in all basins and the time series data for streamflow and evapotranspiration  
 434 predictions and observations of the STL and MTL models in the basin exhibiting the most  
 435 considerable improvement in streamflow prediction. The 1:1 line is represented as a black  
 436 dashed line in both (a) and (b), where points above the line denote a higher NSE using the MTL  
 437 model compared to the STL model. In (c) and (d), the evapotranspiration observational data is  
 438 showcased as a scatterplot, with observations gathered at eight-day intervals, while the  
 439 streamflow observational data and predicted values for streamflow and evapotranspiration  
 440 variables are provided daily.

441 In conclusion, the prediction performance of MTL and STL models was generally  
 442 comparable. In most instances, the prediction of each variable was either nearly equal to or  
 443 slightly superior to that of the STL models. Under certain loss weight configurations, the MTL  
 444 model might exhibit marginally superior performance. Furthermore, instead of training and  
 445 deploying multiple STL models, it was more efficient to select a unified loss weight and utilize a  
 446 single MTL model to simulate the multi-variate hydrological process and predict multiple  
 447 outputs concurrently.

### 448 3.2 Temporal and spatial generalization of MTL and STL models

449 We extended the comparison of STL and MTL models to examine how temporal  
 450 generalizability evolved with an increased number of trained basins and assessed spatial  
 451 generalizability through a PUB test. Figure 5 depicts the scaling curves of both models. Due to  
 452 the spatial extrapolation, blue lines in Figures 5 generally display lower NSE values than red  
 453 lines.



455  
 456 **Figure 5.** Scaling curves for MTL and STL models. The STL models are indicated by dash lines,  
 457 while the MTL models are represented by solid lines. Temporal testing results are represented by

458 red lines and spatial testing results are shown in blue lines. Figures (a) show the predictions of  
459 streamflow, while Figure(b) outline the predictions of evapotranspiration. The y-axis signifies  
460 the median NSE, reflecting the mean value of median NSEs across all folds in a particular  
461 setting. The x-axis represents the percentage of basins used for model training in each setting.  
462 We established 11 scenarios for training, which encompass 5%, 10%, 20%, 25%, 33%, 50%,  
463 66%, 75%, 80%, 90% and 95% of the basins.

464 A consistent trend observed across all subplots was the enhancement in median NSE  
465 value with the rising percentage of basins utilized for training, visible in both temporal and  
466 spatial generalization tests. This pattern indicated that an augmented dataset enhanced the  
467 generalization prowess of DL models in hydrological contexts. Fang et al. (2022) linked this  
468 phenomenon to a data synergy effect, suggesting that accumulating and training more  
469 heterogeneous data enabled DL models to generate better predictions. Our spatial generalization  
470 test affirmed this, highlighting that even in a PUB scenario, the diversity of basins could enhance  
471 the prediction accuracy for all hydrological outputs.

472 Moreover, it became evident that spatial generalizability appeared to improve more  
473 markedly than temporal generalizability. For example, considering streamflow (Q) in Figure  
474 5(a), the median NSE values ranged from approximately 0.58 to 0.70 as the basin training  
475 percentage progressed from 5% to 95%. In contrast, in PUB contexts, this range expanded from  
476 about 0.30 to roughly 0.60, reflecting an enhancement of almost 100%. Similar trends are  
477 evident for evapotranspiration (ET) predictions. These observations suggested that  
478 heterogeneous data provides more substantial benefits for PUB, whereas local data is generally  
479 sufficient for local predictions.

480 In comparison to STL models, MTL models exhibited varying performances in temporal  
481 and spatial generalization tests, predominantly in three distinct patterns. In one scenario, for both  
482 streamflow and evapotranspiration, MTL models marginally underperformed. For example, a 2-  
483 fold cross-validation (equivalent to training with 50% of the basins) assessing the PUB  
484 performance of STL and MTL could prematurely suggest that MTL has weaker spatial  
485 generalization capabilities than STL. In another scenario, a trade-off between streamflow and  
486 evapotranspiration resulted in MTL outperforming STL for one variable while underperforming  
487 for the other, as observed in the 80% training data scenario. Yet, there were instances where  
488 MTL models showcased superior prediction, such as when 33% of basins were used for training,  
489 outperforming in both variable predictions. Considering we only chose one ratio for MTL's

490 different task loss weight, there should be some randomness in the results, but after comparing  
491 MTL with STL in these different scenarios, it could be inferred that MTL model won't be worse  
492 in both temporal and spatial generalization than multiple STL models.

### 493 3.3 Reliability assessment via analysis of internal states in MTL and STL models

494 The cell state of the LSTM model serves as a vital tool for retaining and transmitting  
495 information throughout time series, encapsulating the long-term dependencies observed in  
496 sequences. This characteristic facilitates a deeper understanding of the learning process in  
497 hydrological phenomena (Lees et al., 2022). Before diving into the probe analysis, it's helpful to  
498 examine the direct correlation of the internal states with outputs. Figure 6 offers an illustrative  
499 representation of the correlation coefficients between the hidden layer cell states of the LSTM in  
500 the MTL model compared to two STL models, set against observed evapotranspiration data. The  
501 LSTM hidden layer comprises 256 cell state units, across 591 basins.

502 The most prominent correlation between the LSTM cell state and evapotranspiration is  
503 observed in the STL model for evapotranspiration (STL-ET), followed closely by the MTL  
504 model. Conversely, the STL model for streamflow (STL-Q), which excludes evapotranspiration  
505 from its output, exhibits the least correlation. Figure 6(a) indicates that many basins, especially  
506 around cell numbers 0, 50, and 100, have correlation coefficients approaching 1 or -1.  
507 Meanwhile, Figure 6(b) suggests that the MTL model's LSTM cell state maintains a strong  
508 correlation with evapotranspiration, albeit marginally weaker than the STL-ET model. This  
509 difference can be attributed to the shared layer structure of the MTL model compared to the  
510 specialized nature of the STL-ET model. Figure 6(c) emphasizes the subdued correlation  
511 between the STL-Q model's LSTM cell state and evapotranspiration. However, specific cells,  
512 such as those between 128 to 132, display discernible correlation patterns.

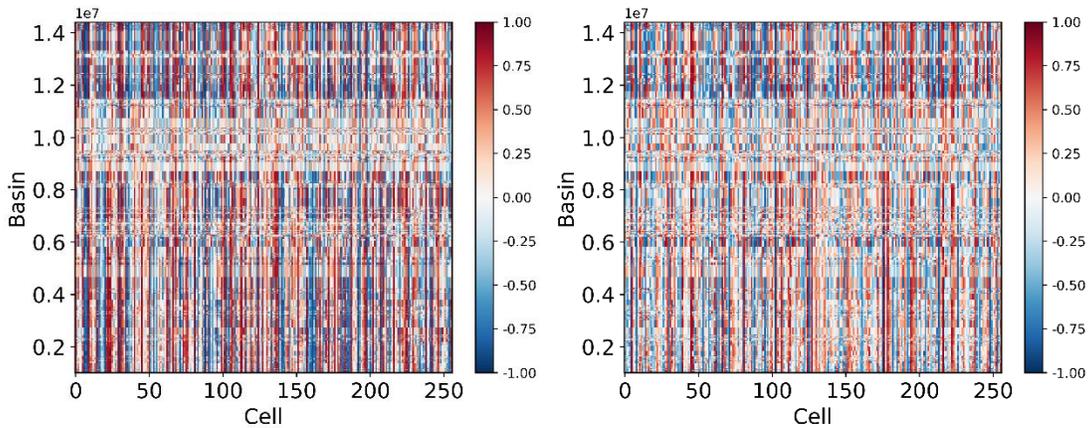
513 We also calculated the maximum absolute value of correlation between cell state and  
514 observation data of evapotranspiration and get the median value of the maximum values for all  
515 basins (median-max-corr). It showed that the values of STL-ET, MTL and STL-Q models are  
516 0.93, 0.93 and 0.85, respectively. Corresponding analyses for the relationship between LSTM  
517 cell states and streamflow are presented in Figure S2 in Supporting Information, where the  
518 median-max-corr values for streamflow in STL-ET, MTL, and STL-Q models are 0.50, 0.71, and  
519 0.71, respectively.

520 The shared-layer LSTM in the MTL model adeptly captured the intricacies of both  
 521 streamflow and evapotranspiration. Although its correlations with individual variables might not  
 522 match those of specialized STL models, its multi-variable proficiency was commendable. STL  
 523 models inherently focus on singular variables, but due to the interrelated nature of hydrological  
 524 components, they might inadvertently capture patterns from non-target variables. In contrast, an  
 525 MTL model, trained on multiple variables, offered a well-rounded correlation pattern with each.  
 526 Simply put, while individual models might discern patterns of related variables due to inherent  
 527 hydrological links, models tailored for multi-variable predictions better comprehend the complex  
 528 interrelationships, even if their correlations appear slightly less intense than singular-focused  
 529 models.

530

(a) STL-ET

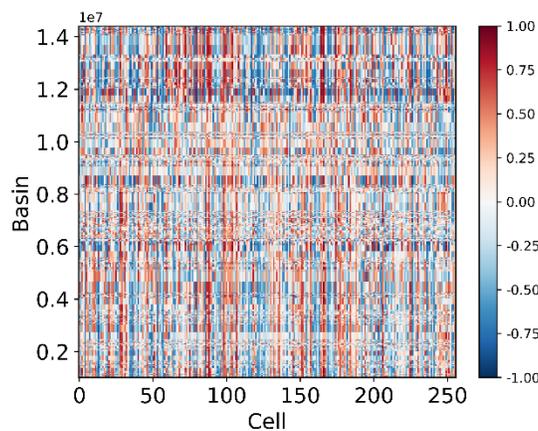
(b) MTL



531

532

(c) STL-Q

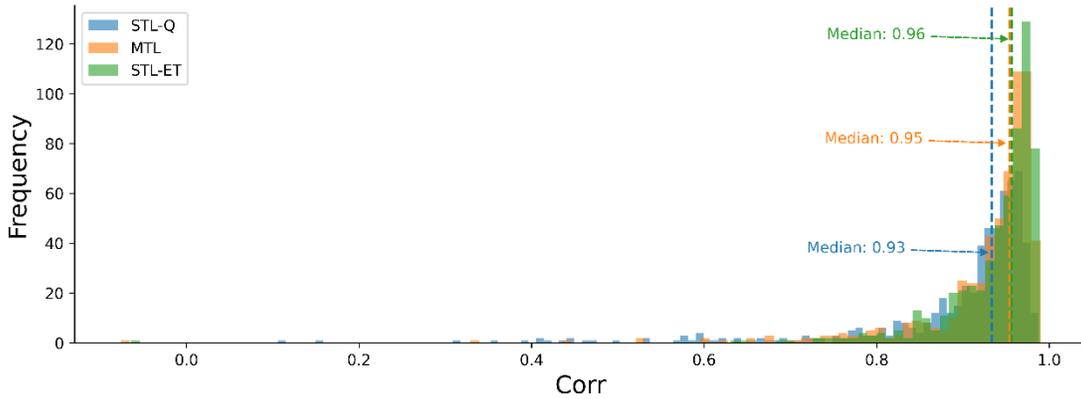


533

534 **Figure 6.** Correlations between the trained LSTM's cell states during the testing period and  
 535 evapotranspiration in different models for each basin. Panels (a), (b), and (c) correspond to the

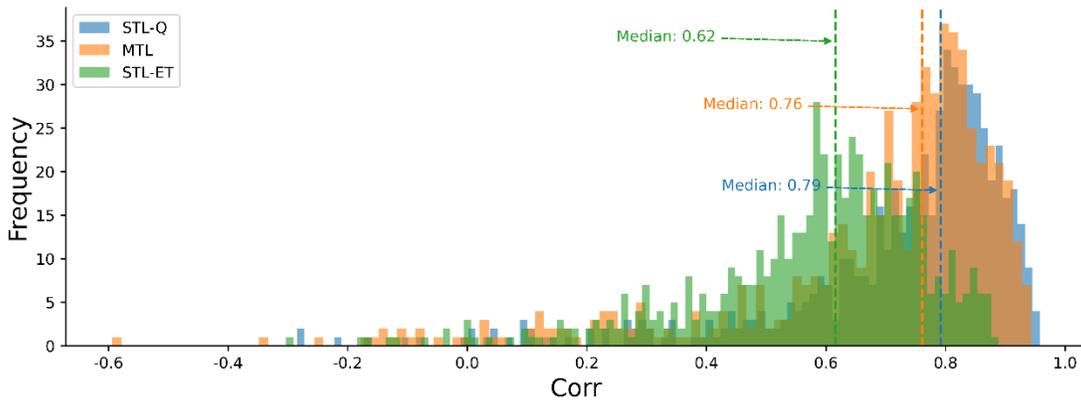
536 STL-ET, MTL, and STL-Q models, respectively. Basins on the y-axis are identified by their 8-  
 537 digit ID from the CAMELS dataset, where notation such as "1e7" represents "10<sup>7</sup>", and "0.2"  
 538 corresponds to "02000000". The x-axis cell labels represent the index of the cell unit within the  
 539 LSTM's cell state.

540 (a) Corr of evapotranspiration probe's prediction



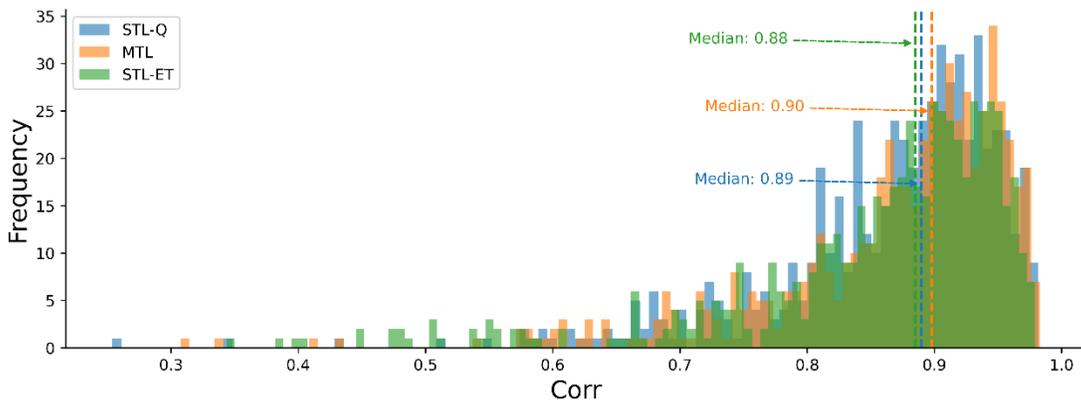
541

542 (b) Corr of streamflow probe's prediction



543

544 (c) Corr of surface soil moisture probe's prediction



545

546 **Figure 7.** A comparison of the correlation coefficients of (a) evapotranspiration (ET), (b)  
547 streamflow (Q), and (c) surface soil moisture (SSM) probes across different models. The blue,  
548 orange, and green bars respectively represent STL-Q, MTL, and STL-ET models. Given that the  
549 correlation serves as a performance indicator, it assumes only positive values, unlike the  
550 correlation between cell states and target variables, which can take negative values.

551 Figures 7(a) and 7(b) illustrate histograms of prediction correlation coefficients for  
552 evapotranspiration and streamflow probes, respectively, across three DL models: STL-Q, MTL,  
553 and STL-ET. The median value of 591 basins for the evapotranspiration probe are approximately  
554 0.93, 0.95, and 0.96 for STL-Q, MTL, and STL-ET models, respectively. For the streamflow  
555 probe, these values are approximately 0.79, 0.76, and 0.62 across the respective models.  
556 Evidently, from the perspective of probe prediction, the LSTM cell state of the STL-ET/STL-Q  
557 model exhibits the strongest correlation with evapotranspiration/streamflow, followed by the  
558 MTL model. In contrast, the LSTM cell state of the STL-Q/STL-ET model shows the weakest  
559 correlation with evapotranspiration/streamflow.

560 The use of cell states to predict a non-target variable was significantly less effective.  
561 While Lees et al. (2022) referred to the performance as “Hydrological Concept Formation” of  
562 DL models and deemed it acceptable, in the absence of constraints imposed by multiple outputs,  
563 the probe’s performance of the STL model might not match that of the MTL model. One  
564 interesting phenomenon was that the highest correlation did not originate from the MTL model.  
565 A linear probe finding the correlation between all cell states and the probe's target variable did  
566 not equate to the highest correlation from one cell state.

567 Figure 7(c) illustrates the correlation coefficients between the predicted and observed  
568 values obtained from the surface soil moisture probe. The MTL model achieves the highest  
569 correlation coefficient for probe prediction results, with a median value for all basins  
570 approximating 0.90. In contrast, the STL-Q and STL-ET models yield correlation coefficients of  
571 approximately 0.89 and 0.88, respectively. This suggested that the shared LSTM layer in the  
572 MTL model, by factoring in input-output correlations for multiple variables, could effectively  
573 learn hydrological processes relevant to non-target variables. Combining all these results, we  
574 proposed that this layer should not simply be considered a trade-off mechanism for multiple  
575 variables. Instead, the learned correlations were more closely aligned with hydrological  
576 processes, thereby enhancing the reliability of the MTL model.

577           Sections 3.1 and 3.2 highlighted that the MTL model, when predicting multiple variables,  
578 was not inferior to the two STL models with large datasets. In fact, it might slightly surpass them  
579 under specific loss weight ratios. Moreover, the MTL model can output multiple variables  
580 concurrently, whereas multiple models would need to be constructed for STL. The results from  
581 Figures 6 and 7 indicated that while the shared-layer LSTM in the MTL model effectively  
582 learned patterns from multiple variables, its individual correlation with each variable might not  
583 be as strong as the dedicated focus each STL model has on its specific target variable. From  
584 these findings, we inferred that the shared LSTM layer in the MTL model excelled at discerning  
585 input-output relationships across multiple variables and delving into variable-specific input-  
586 output correlations. This strengthened the reliability of the correlation rules, a phenomenon we  
587 termed the 'variable synergy effect' within the MTL model framework.

#### 588 **4 Discussion**

589           The "variable synergy" effect, inherent to the MTL model, goes beyond just predicting  
590 multiple targets within a single framework. Generally, MTL models show generalization  
591 capabilities comparable to those achieved by using multiple STL models. The combined layers in  
592 the multi-task neural network often yield a more reliable internal representation compared to the  
593 STL models. The implications of this synergy effect could be interpreted in some way and found  
594 resemblance with the principle of multi-objective optimization (MOP). Such an effect also held  
595 the promise to advance hydrological modeling. We would further explore potential  
596 improvements from MTL models and address the limitations of this study in the subsequent  
597 section.

##### 598           4.1 Trade-off or synergy with multiple outputs in MTL

599           Both STL and MTL models employ deep learning as a universal approximator to capture  
600 the intricate, high-dimensional relationships between inputs and outputs. However, MTL  
601 distinguishes itself by utilizing the relationships between multiple outputs. Through assimilating  
602 loss from these interrelated outputs and updating the shared layers of the neural network, MTL  
603 can aggregate and leverage shared information across tasks. Then, the internal states of an MTL  
604 neural network show a stronger correlation with third-party water balance components,

605 indicating a more comprehensive representation of basin hydrological processes compared to  
606 STL models.

607 In certain MTL studies, the input data itself (Le et al., 2018) or specific noise within the  
608 input data (Pironkov et al., 2017) can serve as auxiliary tasks, effectively acting as regularization  
609 methods. These can improve the generalization performance of the main target prediction.  
610 Hence, in MTL modeling, the inclusion of other tasks can be seen as a regularization method that  
611 reduces overfitting.

612 According to results in Figure 4 in section 3.1, we could find that in MTL modeling, the  
613 notion of trade-offs is salient; bolstered prediction performance for one variable might entail  
614 minor setbacks for another. This dynamic resembles the trade-offs seen in MOP, typical of  
615 reservoir operations. MTL inherently involves an MOP process, and the strategy employed in  
616 this study for MOP involves using weights to convert multi-target objectives into single-target  
617 objectives. A more nuanced or flexible approach could involve strategies like NSGA-II to create  
618 a Pareto front, providing a clear visualization of the competitive relationship between different  
619 targets. Yet, for about  $10^5$  parameters, traditional evolutionary algorithms fall short. This  
620 indicates a prospect for investigating MOP-adapted stochastic gradient descent algorithms in  
621 upcoming studies.

622 Interestingly, evapotranspiration (ET) and streamflow (Q) are not as conflicting as water  
623 supply and flood control in reservoir operations (Castelletti et al., 2013). This insight offers deep  
624 learning researchers a unique lens. They might probe deeper, analyzing and gleaning input data  
625 for latent variables from extant hydrological process insights. Gathering these resources could  
626 enable a Pareto improvement for the multi-variable learning process, potentially enhancing  
627 overall model performance.

#### 628 4.2 Potential and limitations of MTL models

629 MTL models prove advantageous in modeling variables that incur significant  
630 observational costs, especially when paired with more affordably observed variables within one  
631 model. For instance, Surface Soil Moisture (SSM) typically has a shorter observation period  
632 compared to streamflow, which prompts the exploration of integrating the longer streamflow  
633 data into the multi-task learning model to enhance data effectiveness. This approach explores the

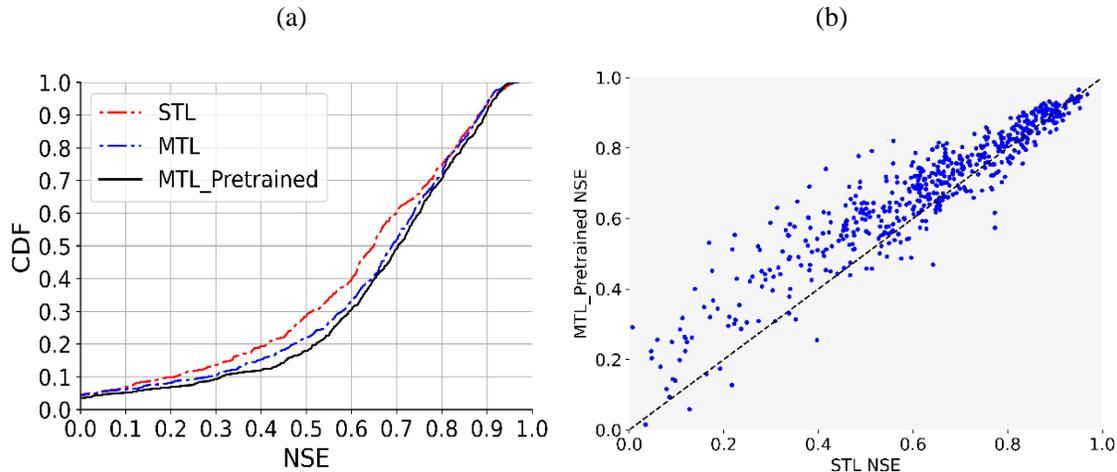
634 leverage of long-term data to forecast short-term data within the MTL paradigm, termed as  
635 "data-augmentation with variable synergy". Evapotranspiration was not included in this  
636 exploration because, based on our preliminary analysis, its predictive performance was already  
637 outstanding, leaving minimal room for enhancement through this technique.

638         Initially, we pretrained the MTL model using only the streamflow data from 2005-04-01  
639 to 2015-03-31, a period without any SSM data records. This treated the MTL model as an STL  
640 model. To ensure the model focused solely on the streamflow, the non-shared fully connected  
641 layer dedicated to the SSM task was intentionally ignored, and its associated loss weight was set  
642 to zero. Subsequently, the MTL model underwent further training using both streamflow and  
643 SSM data from 2015-04-01 to 2020-09-30. We termed model with this training strategy as  
644 MTL\_Pretrained, as depicted in Figure 8. It was then compared to a standard MTL model, which  
645 was trained on both streamflow and SSM data over the same period without any pretraining, as  
646 well as the STL model for SSM. The outcomes of these comparisons are presented in Figure 8.

647         As Figure 8a illustrates, modeling SSM over short durations with limited datasets poses  
648 challenges. However, using an MTL modeling framework can significantly improve the  
649 prediction of SSM, utilizing the synergy effect from streamflow and SSM. Through pretraining,  
650 the LSTM weights and bias are first calibrated guided by streamflow data, circumventing the  
651 commencement from entirely random states. This pretrained model, when retrained, could lead  
652 to slightly better prediction performance. Therefore, even if there is no observation for the data-  
653 scarce variable, it is recommended to use a trained model for another data-rich variable as the  
654 pretrained model, rather than random initialization of bias and weights. Figure 8b confirms that  
655 the pretrained MTL model outperforms in the majority of basins.

656

657



658

659 **Figure 8.** Demonstrating the augment effect for the data-scarce variable from the data-rich  
 660 variable within the MTL modeling framework. Figure 8a is the empirical cumulative density  
 661 function plot for three models: an STL of SSM (STL), an MTL for SSM and Q (MTL) and the  
 662 MTL\_Pretrained model. Figure 8b demonstrates the comparison of NSEs between the STL  
 663 model and the MTL\_Pretrained model. The black line represents a 1:1 line. Points above this line  
 664 indicate that the MTL\_Pretrained NSE is superior.

665 This study proposes an empirical rule that an increased number of observed variables can  
 666 potentially enhance the prediction of less-observed variables within an MTL model. This finding  
 667 is particularly beneficial for predicting hydrological variables with fewer observations, such as  
 668 groundwater streamflow. Under this circumstance, for the high-cost observations, we could use  
 669 more weak-labeled data such as crowdsource data, they can be involved in the multi-task  
 670 modeling framework and provide more information to calibrate the model, which is very difficult  
 671 in traditional modeling methods.

672 Although multiple output variables can bring about predictive refinements, realizing  
 673 substantial advancements without additional input remains a challenge. Sadler et al. (2022)  
 674 suggested that optimizing the loss weight for each variable on a basin-specific basis could further  
 675 improve the prediction within the MTL model framework, but the observed improvements were  
 676 still not significant. This limitation stems from reaching a local optimal point in the feasible  
 677 region of the high-dimensional parameter space without additional information. Since the  
 678 predictions did not improve considerably, the impact on the long-term water balance was minor,  
 679 even though more water balance components were included in the outputs.

680 In hydrological deep learning models, it becomes pivotal to integrate deeper insights into  
 681 the rainfall-runoff dynamics, like pre-event soil moisture. By integrating this additional

682 information, we can gain a more comprehensive understanding of hydrological processes, which  
683 can, in turn, improve the accuracy of predictions.

## 684 **5 Conclusions**

685 This paper explored the role of multi-task deep learning in hydrological modeling across  
686 591 catchments in the CAMELS dataset, using remote sensing observations of actual  
687 evapotranspiration and ground-based streamflow data. An MTL model, rooted in the LSTM  
688 neural network architecture, was developed. We evaluated each variable's predictive  
689 performance of the MTL model by contrasting it with those of two STL models in terms of both  
690 temporal and spatial generalizability. The correlation coefficients between the LSTM cell states  
691 of each model and their corresponding output variables were further investigated. Then a surface  
692 soil moisture probe, which enabled an examination of the neural network's ability to extract  
693 internal representations for the hydrological process was also constructed.

694 Our findings demonstrate that the MTL model, designed for simultaneous predictions of  
695 multiple outputs, consistently matched the performance metrics of its STL counterparts. In  
696 contrast, STL models are restricted to predicting a single output variable, limiting their ability to  
697 capture associations between hydrologic variables. Moreover, in both temporal and spatial  
698 generalization contexts, the MTL model exhibited performance comparable with STL models,  
699 regardless of the dataset size. This highlights the robustness of MTL within hydrological  
700 modeling frameworks, underscoring the resilience of multi-task learning in hydrological  
701 modeling. As a result, the MTL model emerges as a promising deep learning instrument for  
702 further hydrological process exploration, and may soon become the preferred approach in  
703 hydrological modeling over STL.

704 Regarding model reliability, the MTL model mines the relevance of multiple variables  
705 without a marked bias towards any single target, unlike the STL model. Though the MTL's  
706 shared-layer LSTM might have a marginally reduced correlation for individual variables  
707 compared to STL models, it still upholds a reasonable correlation with observations for various  
708 variables. On the other hand, the STL model's correlation with non-target variable observations  
709 is notably weaker. Additionally, the LSTM cell states of the MTL model align more closely with  
710 hydrological processes than those of the STL models. A probe designed for SSM using LSTM  
711 cell states—excluded from all model training—highlighted a superior prediction correlation in

712 the MTL model. This suggests that MTL models better bridge inputs with multiple outputs,  
713 while STL models concentrate mainly on specific target variables.

714 The MTL model also showcased its potential as a regular deep learning method,  
715 especially when faced with limited data observations for certain variables. Its adaptability is  
716 particularly beneficial for bridging data gaps. Nevertheless, a deeper exploration into the  
717 connection between MTL and multi-objective optimization is required. Leveraging gradient-  
718 based multi-objective optimization methods to identify the Pareto frontier could push the  
719 frontiers of MTL in hydrological modeling. Another critical challenge for deep learning in  
720 hydrology remains the need for comprehensive data. It's crucial to understand that hydrological  
721 processes extend beyond just meteorological influences. Incorporating a wider range of ground-  
722 based hydrological time-series data, including pre-event soil moisture, can refine our  
723 understanding of hydrological patterns, driving more precise predictions. In summary, the future  
724 of hydrological modeling will benefit from blending deep learning with multi-objective  
725 optimization techniques, leveraging vast and diverse datasets for richer insights.

## 726 **Acknowledgments**

727 W. Ouyang and L. Ye were primarily supported by the General Program (Grant No.  
728 5217090835) from the National Natural Science Foundation of China (NSFC). W. Ouyang was  
729 also supported by the Youth Science Foundation Project (Grant No. 52309010) from NSFC. All  
730 authors contributed to writing-reviewing and editing, with primary contributions by W. Ouyang.  
731 The contact author has declared that none of the authors has any competing interests.

732

## 733 **Data Availability Statement**

734 All data used in this study are available from public sources. The NLDAS-II dataset can  
735 be downloaded from the website (<http://dx.doi.org/10.5067/THUF4J1RLSYG>), which originally  
736 obtained the dataset from the NOAA/NCEP; The basin attribute data can be downloaded from  
737 the CAMELS website (<http://dx.doi.org/10.5065/D6G73C3Q>) provided by the U.S. National  
738 Center for Atmospheric Research; The SMAP surface soil moisture dataset is available at  
739 (<https://doi.org/10.5067/ZX7YX2Y2LHEB>); The streamflow data can be obtained from USGS

740 Water Data for the Nation website (<http://dx.doi.org/10.5066/F7P55KJN>). The code used in this  
741 study are available in the open-source repository (<https://doi.org/10.5281/zenodo.10024012>)

742

## 743 **References**

744 Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017). The CAMELS data set:  
745 Catchment attributes and meteorology for large-sample studies. *Hydrology and Earth  
746 System Sciences*, *21*(10), 5293–5313. <https://doi.org/10.5194/hess-21-5293-2017>

747 Bannai, T., Xu, H., Utsumi, N., Koo, E., Lu, K., & Kim, H. (2023). Multi-Task Learning for  
748 Simultaneous Retrievals of Passive Microwave Precipitation Estimates and Rain/No-Rain  
749 Classification. *Geophysical Research Letters*, *50*(7), e2022GL102283.  
750 <https://doi.org/10.1029/2022GL102283>

751 Becker, R., Koppa, A., Schulz, S., Usman, M., Beek, T. aus der, & Schüth, C. (2019). Spatially  
752 distributed model calibration of a highly managed hydrological system using remote  
753 sensing-derived ET data. *Journal of Hydrology*, *577*, 123944.  
754 <https://doi.org/10.1016/j.jhydrol.2019.123944>

755 Belinkov, Y., Durrani, N., Dalvi, F., Sajjad, H., & Glass, J. (2017). What do Neural Machine  
756 Translation Models Learn about Morphology? *Proceedings of the 55th Annual Meeting of  
757 the Association for Computational Linguistics (Volume 1: Long Papers)*, 861–872.  
758 <https://doi.org/10.18653/v1/P17-1080>

759 Castelletti, A., Pianosi, F., & Restelli, M. (2013). A multiobjective reinforcement learning  
760 approach to water resources systems operation: Pareto frontier approximation in a single  
761 run. *Water Resources Research*, *49*(6), 3476–3486. <https://doi.org/10.1002/wrcr.20295>

- 762 Cipolla, R., Gal, Y., & Kendall, A. (2018). Multi-task Learning Using Uncertainty to Weigh  
763 Losses for Scene Geometry and Semantics. *2018 IEEE/CVF Conference on Computer  
764 Vision and Pattern Recognition*, 7482–7491. <https://doi.org/10.1109/CVPR.2018.00781>
- 765 Dembélé, M., Ceperley, N., Zwart, S. J., Salvadore, E., Mariethoz, G., & Schaepli, B. (2020).  
766 Potential of satellite and reanalysis evaporation datasets for hydrological modelling under  
767 various model calibration strategies. *Advances in Water Resources*, *143*, 103667.  
768 <https://doi.org/10.1016/j.advwatres.2020.103667>
- 769 Dembélé, M., Hrachowitz, M., Savenije, H. H. G., Mariéthoz, G., & Schaepli, B. (2020).  
770 Improving the Predictive Skill of a Distributed Hydrological Model by Calibration on  
771 Spatial Patterns With Multiple Satellite Data Sets. *Water Resources Research*, *56*(1).  
772 <https://doi.org/10.1029/2019WR026085>
- 773 Fang, K., Kifer, D., Lawson, K., Feng, D., & Shen, C. (2022). The Data Synergy Effects of  
774 Time-Series Deep Learning Models in Hydrology. *Water Resources Research*, *58*(4),  
775 e2021WR029583. <https://doi.org/10.1029/2021WR029583>
- 776 Feng, D., Fang, K., & Shen, C. (2020). Enhancing Streamflow Forecast and Extracting Insights  
777 Using Long-Short Term Memory Networks With Data Integration at Continental Scales.  
778 *Water Resources Research*, *56*(9), e2019WR026793.  
779 <https://doi.org/10.1029/2019WR026793>
- 780 Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., & Moore, R. (2017). Google  
781 Earth Engine: Planetary-scale geospatial analysis for everyone. *Big Remotely Sensed  
782 Data: Tools, Applications and Experiences*, *202*, 18–27.  
783 <https://doi.org/10.1016/j.rse.2017.06.031>

- 784 Guo, M., Haque, A., Huang, D.-A., Yeung, S., & Fei-Fei, L. (2018, September). Dynamic Task  
785 Prioritization for Multitask Learning. *Proceedings of the European Conference on*  
786 *Computer Vision (ECCV)*.
- 787 Herman, M. R., Nejadhashemi, A. P., Abouali, M., Hernandez-Suarez, J. S., Daneshvar, F.,  
788 Zhang, Z., Anderson, M. C., Sadeghi, A. M., Hain, C. R., & Sharifi, A. (2018).  
789 Evaluating the role of evapotranspiration remote sensing data in improving hydrological  
790 modeling predictability. *Journal of Hydrology*, 556, 39–49.  
791 <https://doi.org/10.1016/j.jhydrol.2017.11.009>
- 792 Hewitt, J., & Liang, P. (2019). Designing and Interpreting Probes with Control Tasks.  
793 *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*  
794 *Processing and the 9th International Joint Conference on Natural Language Processing*  
795 *(EMNLP-IJCNLP)*, 2733–2743. <https://doi.org/10.18653/v1/D19-1275>
- 796 Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8),  
797 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- 798 Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are  
799 universal approximators. *Neural Networks*, 2(5), 359–366. <https://doi.org/10.1016/0893->  
800 [6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8)
- 801 Hu, X., Shi, L., Lin, G., & Lin, L. (2021). Comparison of physical-based, data-driven and hybrid  
802 modeling approaches for evapotranspiration estimation. *Journal of Hydrology*, 601,  
803 126592. <https://doi.org/10.1016/j.jhydrol.2021.126592>
- 804 Jiang, S., Zheng, Y., & Solomatine, D. (2020). Improving AI System Awareness of Geoscience  
805 Knowledge: Symbiotic Integration of Physical Approaches and Deep Learning.

- 806 *Geophysical Research Letters*, 47(13), e2020GL088229.  
807 <https://doi.org/10.1029/2020GL088229>
- 808 Kraft, B., Jung, M., Körner, M., & Reichstein, M. (2020). Hybrid modeling: Fusion of a deep  
809 learning approach and a physics-based model for global hydrological modeling. *The*  
810 *International Archives of the Photogrammetry, Remote Sensing and Spatial Information*  
811 *Sciences*, XLIII-B2-2020, 1537–1544. [https://doi.org/10.5194/isprs-archives-XLIII-B2-](https://doi.org/10.5194/isprs-archives-XLIII-B2-2020-1537-2020)  
812 [2020-1537-2020](https://doi.org/10.5194/isprs-archives-XLIII-B2-2020-1537-2020)
- 813 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019).  
814 Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine  
815 Learning. *Water Resources Research*, 55(12), 11344–11354.  
816 <https://doi.org/10.1029/2019WR026065>
- 817 Kratzert, F., Klotz, D., Hochreiter, S., & Nearing, G. S. (2021). A note on leveraging synergy in  
818 multiple meteorological data sets with deep learning for rainfall–runoff modeling.  
819 *Hydrology and Earth System Sciences*, 25(5), 2685–2703. [https://doi.org/10.5194/hess-](https://doi.org/10.5194/hess-25-2685-2021)  
820 [25-2685-2021](https://doi.org/10.5194/hess-25-2685-2021)
- 821 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019). Towards  
822 learning universal, regional, and local hydrological behaviors via machine learning  
823 applied to large-sample datasets. *Hydrology and Earth System Sciences*, 23(12), 5089–  
824 5110. <https://doi.org/10.5194/hess-23-5089-2019>
- 825 Le, L., Patterson, A., & White, M. (2018). Supervised autoencoders: Improving generalization  
826 performance with unsupervised regularizers. In S. Bengio, H. Wallach, H. Larochelle, K.  
827 Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in Neural Information*

- 828 *Processing Systems* (Vol. 31). Curran Associates, Inc.  
829 [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/2a38a4a9316c49e5a833517c45](https://proceedings.neurips.cc/paper_files/paper/2018/file/2a38a4a9316c49e5a833517c45)  
830 [d31070-Paper.pdf](#)
- 831 Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve, P.,  
832 Slater, L., & Dadson, S. J. (2022). Hydrological concept formation inside long short-term  
833 memory (LSTM) networks. *Hydrology and Earth System Sciences*, 26(12), 3079–3101.  
834 <https://doi.org/10.5194/hess-26-3079-2022>
- 835 Li, B., Li, R., Sun, T., Gong, A., Tian, F., Khan, M. Y. A., & Ni, G. (2023). Improving LSTM  
836 hydrological modeling with spatiotemporal deep learning and multi-task learning: A case  
837 study of three mountainous areas on the Tibetan Plateau. *Journal of Hydrology*, 620,  
838 129401. <https://doi.org/10.1016/j.jhydrol.2023.129401>
- 839 Li, X., Khandelwal, A., Jia, X., Cutler, K., Ghosh, R., Renganathan, A., Xu, S., Tayal, K., Nieber,  
840 J., Duffy, C., Steinbach, M., & Kumar, V. (2022). Regionalization in a Global Hydrologic  
841 Deep Learning Model: From Physical Descriptors to Random Vectors. *Water Resources*  
842 *Research*, 58(8), e2021WR031794. <https://doi.org/10.1029/2021WR031794>
- 843 Liu, J., Hughes, D., Rahmani, F., Lawson, K., & Shen, C. (2023). Evaluating a global soil  
844 moisture dataset from a multitask model (GSM3 v1.0) with potential applications for  
845 crop threats. *Geoscientific Model Development*, 16(5), 1553–1567.  
846 <https://doi.org/10.5194/gmd-16-1553-2023>
- 847 Liu, J., Rahmani, F., Lawson, K., & Shen, C. (2022). A Multiscale Deep Learning Model for Soil  
848 Moisture Integrating Satellite and In Situ Data. *Geophysical Research Letters*, 49(7),  
849 e2021GL096847. <https://doi.org/10.1029/2021GL096847>

- 850 Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., Sharma, A., & Shen, C. (2021).  
851 Transferring Hydrologic Data Across Continents – Leveraging Data-Rich Regions to  
852 Improve Hydrologic Prediction in Data-Sparse Regions. *Water Resources Research*,  
853 57(5), e2020WR028600. <https://doi.org/10.1029/2020WR028600>
- 854 McCabe, M. F., Rodell, M., Alsdorf, D. E., Miralles, D. G., Uijlenhoet, R., Wagner, W., Lucieer,  
855 A., Houborg, R., Verhoest, N. E. C., Franz, T. E., Shi, J., Gao, H., & Wood, E. F. (2017).  
856 The future of Earth observation in hydrology. *Hydrology and Earth System Sciences*,  
857 21(7), 3879–3914. <https://doi.org/10.5194/hess-21-3879-2017>
- 858 Mladenova, I. E., Bolten, J. D., Crow, W. T., Sazib, N., Cosh, M. H., Tucker, C. J., & Reynolds,  
859 C. (2019). Evaluating the Operational Application of SMAP for Global Agricultural  
860 Drought Monitoring. *IEEE Journal of Selected Topics in Applied Earth Observations and*  
861 *Remote Sensing*, 12(9), 3387–3397. <https://doi.org/10.1109/JSTARS.2019.2923555>
- 862 NASA. (2018). Forcing Data for Phase 2 of the North American Land Data Assimilation System  
863 (NLDAS-2). In *National Aeronautics and Space Administration (NASA)*.  
864 <https://ldas.gsfc.nasa.gov/nldas/NLDAS2forcing.php#AppendixC>
- 865 Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C.,  
866 & Gupta, H. V. (2021). What Role Does Hydrological Science Play in the Age of  
867 Machine Learning? *Water Resources Research*, 57(3), e2020WR028091.  
868 <https://doi.org/10.1029/2020WR028091>
- 869 Omernik, J. M., & Griffith, G. E. (2014). Ecoregions of the conterminous United States:  
870 Evolution of a hierarchical spatial framework. *Environmental Management*, 54(6), 1249–  
871 1266. <https://doi.org/10.1007/s00267-014-0364-1>

- 872 Ouyang, W., Lawson, K., Feng, D., Ye, L., Zhang, C., & Shen, C. (2021). Continental-scale  
873 streamflow modeling of basins with reservoirs: Towards a coherent deep-learning-based  
874 strategy. *Journal of Hydrology*, 599, 126455.  
875 <https://doi.org/10.1016/j.jhydrol.2021.126455>
- 876 Pironkov, G., Dupont, S., Wood, S. U. N., & Dutoit, T. (2017). Noise and Speech Estimation as  
877 Auxiliary Tasks for Robust Speech Recognition. In N. Camelin, Y. Estève, & C. Martín-  
878 Vide (Eds.), *Statistical Language and Speech Processing* (pp. 181–192). Springer  
879 International Publishing. [https://doi.org/10.1007/978-3-319-68456-7\\_15](https://doi.org/10.1007/978-3-319-68456-7_15)
- 880 Rahmani, F., Lawson, K., Ouyang, W., Appling, A., Oliver, S., & Shen, C. (2021). Exploring the  
881 exceptional performance of a deep learning stream temperature model and the value of  
882 streamflow data. *Environmental Research Letters*, 16(2), 024025.  
883 <https://doi.org/10.1088/1748-9326/abd501>
- 884 Running, S., Mu, Q., & Zhao, M. (2017). *MOD16A2 MODIS/Terra Net Evapotranspiration 8-*  
885 *Day L4 Global 500m SIN Grid V006* [dataset]. NASA EOSDIS Land Processes DAAC.  
886 <https://doi.org/10.5067/MODIS/MOD16A2.006>
- 887 Sadler, J. M., Appling, A. P., Read, J. S., Oliver, S. K., Jia, X., Zwart, J. A., & Kumar, V. (2022).  
888 Multi-Task Deep Learning of Daily Streamflow and Water Temperature. *Water Resources*  
889 *Research*, 58(4), e2021WR030138. <https://doi.org/10.1029/2021WR030138>
- 890 Schmidt, L., Heße, F., Attinger, S., & Kumar, R. (2020). Challenges in Applying Machine  
891 Learning Models for Hydrological Inference: A Case Study for Flooding Events Across  
892 Germany. *Water Resources Research*, 56(5), e2019WR025924.  
893 <https://doi.org/10.1029/2019WR025924>

- 894 Sener, O., & Koltun, V. (2018). Multi-Task Learning as Multi-Objective Optimization. In S.  
895 Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.),  
896 *Advances in Neural Information Processing Systems* (Vol. 31). Curran Associates, Inc.  
897 <https://proceedings.neurips.cc/paper/2018/file/432aca3a1e345e339f35a30c8f65edce->  
898 [Paper.pdf](#)
- 899 Shah, S., Duan, Z., Song, X., Li, R., Mao, H., Liu, J., Ma, T., & Wang, M. (2021). Evaluating the  
900 added value of multi-variable calibration of SWAT with remotely sensed  
901 evapotranspiration data for improving hydrological modeling. *Journal of Hydrology*, 603,  
902 127046. <https://doi.org/10.1016/j.jhydrol.2021.127046>
- 903 Shen, C. (2018). A Transdisciplinary Review of Deep Learning Research and Its Relevance for  
904 Water Resources Scientists. *Water Resources Research*, 54(11), 8558–8593.  
905 <https://doi.org/10.1029/2018WR022643>
- 906 Standley, T., Zamir, A., Chen, D., Guibas, L., Malik, J., & Savarese, S. (2020). Which Tasks  
907 Should Be Learned Together in Multi-task Learning? In H. D. III & A. Singh (Eds.),  
908 *Proceedings of the 37th International Conference on Machine Learning* (Vol. 119, pp.  
909 9120–9132). PMLR. <https://proceedings.mlr.press/v119/standley20a.html>
- 910 Tobin, K. J., & Bennett, M. E. (2017). Constraining SWAT Calibration with Remotely Sensed  
911 Evapotranspiration Data. *JAWRA Journal of the American Water Resources Association*,  
912 53(3), 593–604. <https://doi.org/10.1111/1752-1688.12516>
- 913 Tong, R., Parajka, J., Salentinig, A., Pfeil, I., Komma, J., Széles, B., Kubáň, M., Valent, P.,  
914 Vreugdenhil, M., Wagner, W., & Blöschl, G. (2021). The value of ASCAT soil moisture  
915 and MODIS snow cover data for calibrating a conceptual hydrologic model. *Hydrology*

- 916 *and Earth System Sciences*, 25(3), 1389–1410. <https://doi.org/10.5194/hess-25-1389->  
917 2021
- 918 Tong, R., Parajka, J., Széles, B., Greimeister-Pfeil, I., Vreugdenhil, M., Komma, J., Valent, P., &  
919 Blöschl, G. (2022). The value of satellite soil moisture and snow cover data for the  
920 transfer of hydrological model parameters to ungauged sites. *Hydrology and Earth*  
921 *System Sciences*, 26(7), 1779–1799. <https://doi.org/10.5194/hess-26-1779-2022>
- 922 Tsai, W.-P., Feng, D., Pan, M., Beck, H., Lawson, K., Yang, Y., Liu, J., & Shen, C. (2021). From  
923 calibration to parameter learning: Harnessing the scaling effects of big data in  
924 geoscientific modeling. *Nature Communications*, 12(1), 5988.  
925 <https://doi.org/10.1038/s41467-021-26107-z>
- 926 USGS. (2019). National Water Information System: Web interface. In *United States Geological*  
927 *Survey*. <https://waterdata.usgs.gov/nwis/>
- 928 Vandenhende, S., Georgoulis, S., De Brabandere, B., & Van Gool, L. (2019). *Branched Multi-*  
929 *Task Networks: Deciding What Layers To Share* (1904.02920v5). arXiv.  
930 <https://doi.org/10.48550/arXiv.1904.02920>
- 931 Vandenhende, S., Georgoulis, S., Van Gansbeke, W., Proesmans, M., Dai, D., & Van Gool, L.  
932 (2022). Multi-Task Learning for Dense Prediction Tasks: A Survey. *IEEE Transactions on*  
933 *Pattern Analysis and Machine Intelligence*, 44(7), 3614–3633.  
934 <https://doi.org/10.1109/TPAMI.2021.3054719>
- 935 Vu, M. T., & Jardani, A. (2022). Multi-task neural network in hydrological tomography to map  
936 the transmissivity and storativity simultaneously: HT-XNET. *Journal of Hydrology*, 612,  
937 128167. <https://doi.org/10.1016/j.jhydrol.2022.128167>

- 938 Xu, T., Guo, Z., Xia, Y., Ferreira, V. G., Liu, S., Wang, K., Yao, Y., Zhang, X., & Zhao, C.  
939 (2019). Evaluation of twelve evapotranspiration products from machine learning, remote  
940 sensing and land surface models over conterminous United States. *Journal of Hydrology*,  
941 578, 124105. <https://doi.org/10.1016/j.jhydrol.2019.124105>
- 942 Yassin, F., Razavi, S., Wheeler, H., Sapriza-Azuri, G., Davison, B., & Pietroniro, A. (2017).  
943 Enhanced identification of a hydrologic model using streamflow and satellite water  
944 storage data: A multicriteria sensitivity analysis and optimization approach. *Hydrological*  
945 *Processes*, 31(19), 3320–3333. <https://doi.org/10.1002/hyp.11267>
- 946 Yeste, P., Melsen, L. A., & Garc, M. (2023). A Pareto-based sensitivity analysis and multi-  
947 objective calibration approach for integrating streamflow and evaporation data. *Water*  
948 *Resources Research*, 1–29. <https://doi.org/10.1029/2022WR033235>
- 949 Yokoo, K., Ishida, K., Ercan, A., Tu, T., Nagasato, T., Kiyama, M., & Amagasaki, M. (2022).  
950 Capabilities of deep learning models on learning physical relationships: Case of rainfall-  
951 runoff modeling with LSTM. *Science of The Total Environment*, 802, 149876.  
952 <https://doi.org/10.1016/j.scitotenv.2021.149876>
- 953 Zeiler, M. D. (2012). Adadelata: An adaptive learning rate method. *arXiv Preprint*  
954 *arXiv:1212.5701*.
- 955 Zhang, Y., & Yang, Q. (2022). A Survey on Multi-Task Learning. *IEEE Transactions on*  
956 *Knowledge and Data Engineering*, 34(12), 5586–5609.  
957 <https://doi.org/10.1109/TKDE.2021.3070203>
- 958 Zhao, R.-J. (1992). The Xinanjiang model applied in China. *Journal of Hydrology*, 135(1), 371–  
959 381. [https://doi.org/10.1016/0022-1694\(92\)90096-E](https://doi.org/10.1016/0022-1694(92)90096-E)

960 Zhi, W., Ouyang, W., Shen, C., & Li, L. (2023). Temperature outweighs light and flow as the  
961 predominant driver of dissolved oxygen in US rivers. *Nature Water*, 1(3), 249–260.  
962 <https://doi.org/10.1038/s44221-023-00038-z>

963