

Ensembling Differentiable Process-based and Data-driven Models with Diverse Meteorological Forcing Datasets to Advance Streamflow Simulation

Peijun Li¹, Yalan Song¹, Ming Pan², Kathryn Lawson¹, Chaopeng Shen¹

¹Civil and Environmental Engineering, The Pennsylvania State University, University Park, PA, USA

²Center for Western Weather and Water Extremes, Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA

*Correspondence to: Peijun Li, pql5336@psu.edu; Chaopeng Shen, cshen@engr.psu.edu

Abstract

Streamflow simulations ~~via~~produced by different hydrological models ~~have different~~ ~~features~~exhibit distinct characteristics and can provide valuable information ~~after being~~when ensembled. ~~While~~However, few studies have focused on ensembling simulations ~~via~~from models with significant structural differences and evaluating them under both temporal and spatial tests. Here we systematically evaluated and utilized the simulations from two highly different models with great performances: a purely data-driven long short-term memory (LSTM) network and a physics-informed machine learning (“differentiable”) HBV (Hydrologiska Byråns ~~Vattenavdelning~~Vattenbalansavdelning) model (δ HBV). To effectively display the features of the two models, multiple forcing datasets are employed ~~and utilized in~~ ~~two ways~~;. The results show that the simulations of LSTM and δ HBV have distinct features and complement each other well, leading to better Nash-Sutcliffe model efficiency coefficients (NSE) and improved high-flow and low-flow metrics across all spatiotemporal tests, compared to within-class ensembles. Ensembling models trained on a single forcing outperformed a single model using fused forcings, challenging the paradigm of feeding all available data into a single data-driven model. Most notably, δ HBV significantly enhanced spatial interpolation

Formatted: Header

Style Definition: Header: Font: (Asian) SimSun

Formatted: Normal, Line spacing: single

Formatted: Font color: Auto

when incorporated into LSTM, and provided even more prominent benefits for spatial extrapolation where the LSTM-only ensembles degraded significantly, attesting to the value of the structural constraints in δ HBV. These advances set new benchmark records on the well-known CAMELS (Catchment Attributes and Meteorology for Large-sample Studies) hydrological dataset, reaching median NSE values of ~ 0.83 for the temporal test (densely trained scenario), ~ 0.79 for the ungauged basin test (PUB, Prediction in Ungauged Basins), and ~ 0.70 for the ungauged region test (PUR, Prediction in Ungauged Regions). This study advances our understanding of how various model types, each with distinct mechanisms, can be effectively leveraged alongside multi-source datasets across diverse scenarios.

Formatted: Header

Highlights

Formatted: Font color: Auto

- Combining LSTM and δ HBV with diverse forcings sets new accuracy benchmarks
- Ensembling models with one forcing outperforms merging forcings as an input
- δ HBV and LSTM together always increase NSEs, especially spatial generalization
- δ HBV provides valuable spatial constraints in the deterministic ensemble simulations
- δ HBV and LSTM have different error characteristics that can be offset in an ensemble

Formatted: Outline numbered + Level: 1 + Numbering Style: Bullet + Aligned at: 0" + Indent at: 0.29"

Keywords

Formatted: Font color: Auto

Streamflow simulation, differentiable model, deep learning, hybrid modeling, multi-source fusion

1. Introduction

Formatted: Font color: Auto

Streamflow, a critical component of the global hydrosphere, profoundly influences both human society and natural ecosystems (Lins and Slack, 1999). Accurate simulation and prediction of streamflow yield numerous benefits, including improved flood prevention strategies (Brunner et al., 2021). Hydrological models serve as indispensable tools for achieving this objective and can be traditionally categorized into two types: data-driven models (Feng et al., 2020; Kratzert et al., 2018; Liu et al., 2024; Nearing et al., 2024)(Feng et al., 2020; Kratzert et al., 2018; Liu et al., 2024; Nearing et al., 2024) and process-based (or physically-based) models (Newman et al., 2017; Paul et al., 2021). Data-driven models, exemplified by long short-term memory (LSTM) (Feng et al., 2020; Kratzert et al., 2018) and transformer (Liu et al., 2024) networks, excel in learning patterns from multi-source data (Li et al., 2023b, 2024; Liu et al., 2022; Nearing et al., 2024)(Li et al., 2023b, 2024; Liu et al., 2022; Nearing et al., 2024) and generally achieve high performance. However, they often lack interpretability and

may not resolve extreme values very well (~~Li et al., 2020a; Song et al., 2025~~)(~~Li et al., 2020a; Song et al., 2025b~~). Conversely, process-based models, derived deductively from physical laws or conceptualized views of natural systems, offer insights into internal hydrological processes but may exhibit weaker performance due to structural inadequacies (Li et al., 2020a; Zhang et al., 2019).

To combine the benefits and counteract the weaknesses of these two kinds of models, many efforts have been made to incorporate physical constraints and structures into data-driven models to align with fundamental physical principles, such as mass and water balances (Bandai and Ghezzehei, 2021; Wang et al., 2020; Xie et al., 2021). The most seamless integration uses neural networks to provide parameterizations or missing process representations for process-based models (Aboelyazeed et al., 2023; Bindas et al., 2024; Feng et al., 2022; Jiang et al., 2020; Kraft et al., 2022; Rahmani et al., 2023; Song et al., 2024b; Tsai et al., 2021). These differentiable models (Shen et al., 2023) connect (flexible amounts of) prior physical knowledge to neural networks, and have displayed many advantages, including improved computational efficiency and prediction of untrained variables (Tsai et al., 2021), spatial generalization (Feng et al., 2023b), and representation of extremes (~~Song et al., 2025~~)(~~Song et al., 2025b~~). However, it is also unclear whether current differentiable models, e.g., δ HBV, the Hydrologiska Byråns Vattenbalansavdelning (HBV) model implemented within a differentiable framework (~~Feng et al., 2023b; Shen et al., 2023; Song et al., 2025~~)(~~Feng et al., 2023b; Shen et al., 2023; Song et al., 2025b~~), have unique bias characteristics that are associated with the process-based parts of their structures that cannot be reduced once the equations are prescribed.

Orthogonal to such efforts are ensemble simulations (Yu et al., 2024), which combine many members with different biases and uncertainties to mitigate their respective biases in deterministic predictions. Many previous studies have tried ensemble methods to improve

streamflow (Clark et al., 2016; Zounemat-Kermani et al., 2021) based on many factors, like initial conditions (e.g., initial weights and biases in LSTM (Kratzert et al., 2018)), data used for parameterization (Feng et al., 2021), and objective functions (Lin et al., 2024). These studies generally use one model to generate the differences among the ensemble members. Furthermore, some studies (Dion et al., 2021; Solanki et al., 2025) have utilized simulations from multiple different models but are limited to process-based models ~~and resulted, resulting~~ in ensemble simulations that are better than each individual member. Thus far, however, most studies ~~foeushave focused~~ on ~~the~~ simulations from only similar models or model types, and little work has tested an ensemble across the boundary of model types, ~~especiallly~~ particularly between data-driven, process-based, and hybrid models, especially on a large number of samples. Presumably, if each model has its own unique bias, data-driven and process-based models are likely to exhibit greater differences due to their inherently distinct characteristics. It remains unclear whether ensembling across model types should bring benefits to deterministic predictions. Furthermore, grounded in the process-based model, the differentiable process-based hydrological model, such as δ HBV, significantly enhances performance compared to traditional process-based models, while on the other hand introducing greater uncertainty regarding its potential benefits when ensembled. Moreover, previous studies have primarily focused on evaluating ensemble simulations for temporal predictions. However, streamflow simulation under spatial extrapolation scenarios presents greater challenges, and findings from temporal tests may not be directly applicable in this context.

It is known that the performance of any type of hydrologic model heavily depends on the quality of input data, particularly meteorological forcing data (Bell and Moore, 2000; Yao et al., 2020), and other inputs, like the uncertainties of initial conditions, can be mitigated via warming up (Yu et al., 2019). While independent forcing datasets excel in certain aspects, they each carry different error characteristics (Beck et al., 2017; Behnke et al., 2016; Newman et al.,

2019) and accordingly affect the hydrological models in different ways. In order to fully display the different features between LSTM and δ HBV, multiple forcing datasets could be considered. Given the utilization of multiple forcing datasets, one could choose to use data fusion to combine them into a single coherent model input (Kratzert et al., 2021; Sawadekar et al., 2025), or to pass each forcing dataset through a model and then afterwards combine the multiple outputs in an ensemble. It is not clear which approach is more beneficial.

Considering the knowledge gaps discussed above, we sought to answer several research questions:

1. Will a cross-model-type ensemble of LSTM and δ HBV improve deterministic streamflow prediction more than a within-class ensemble?
2. Is it better to use multiple forcings in one model or to ensemble multiple models, each with a different forcing input?
3. Do process-based equations bring unique value to an ensemble, especially in terms of spatial generalizability?

The remainder of this paper is structured as follows: Sect. 2 outlines the hydrological data and models used in this study, as well as the experimental design. Results and discussions are presented in Sect. 3, with conclusions provided in Sect. 4.

2. ~~Material~~Materials and ~~Methodsmethods~~

2.1. ~~The~~ CAMELS hydrologic dataset

The Catchment Attributes and Meteorology for Large-sample Studies (CAMELS) dataset (Addor et al., 2017) is widely employed for hydrological model evaluation and community benchmarking. The CAMELS dataset encompasses 671 basins distributed across the conterminous United States, with basin sizes ranging from 1 to 25,800 km² (median: 335 km²). This standardized and publicly available dataset serves as a benchmark for evaluating various

Formatted: Header

Formatted: Outline numbered + Level: 1 + Numbering Style: 1, 2, 3, ... + Start at: 1 + Alignment: Left + Aligned at: 0.25" + Indent at: 0.5"

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

hydrological models, with LSTM models trained on this dataset often serving as a reference point for comparing other models (Kratzert et al., 2021). CAMELS provides basin-scale data, including streamflow observations and static basin attributes, as well as forcing datasets from three independent sources: Daymet (Thornton et al., 1997), North American Land Data Assimilation System (NLDAS) (Xia et al., 2012), and Maurer (Maurer et al., 2002). Each of the three meteorological forcing datasets operates at a daily temporal resolution, encompassing precipitation, temperature, vapor pressure, and surface radiation variables, with daily temperature extrema of NLDAS and Maurer supplemented from Kratzert et al. (2021). These three meteorological forcing datasets have methodological distinctions in spatial resolution, data generation approaches, and temporal processing (Behnke et al., 2016; Kratzert et al., 2021). Exemplary plots illustrating the differences among the three meteorological forcing datasets are provided in Appendix B. These features can lead to dataset-specific error characteristics and make them valuable for displaying the distinct features of different model types. All model inputs used in this study are detailed in Table C1.

2.2. Long short-term memory

As one kind of deep learning algorithm, long short-term memory (LSTM) (Hochreiter and Schmidhuber, 1997) has unique structures like hidden states and gates activated by the tanh and sigmoid functions (Li et al., 2023a), respectively. These features enable LSTM to excel in streamflow simulation tasks (Feng et al., 2020; Kratzert et al., 2018; Nearing et al., 2024)(Feng et al., 2020; Kratzert et al., 2018; Nearing et al., 2024). In the current benchmark framework, LSTM models are trained using dynamic atmospheric forcings and static basin attributes as inputs, with streamflow as the target output, making it perform well in both temporal and spatial tests (Figure 1a). In this work, for cross-group comparability, we used the LSTM model and its hyperparameters as reported in Kratzert et al. (2021).

2.3. Differentiable HBV model (δ HBV)

The Hydrologiska Byråns Vattenbalansavdelning (HBV) model is a parsimonious bucket-type hydrologic model that simulates various hydrological variables, including snow water equivalent, soil water, groundwater storage, evapotranspiration, quick flow, baseflow, and total streamflow (Aghakouchak and Habib, 2010; Beck et al., 2020; Bergström, 1976, 1992). Recently demonstrated differentiable HBV (δ HBV) model (Feng et al., 2023b; Shen et al., 2023; Song et al., 2024b) incorporates deep neural networks for both regionalized parameterization and missing process representations within a differentiable programming framework that supports “end-to-end” training (Figure 1b). This innovation enables δ HBV to effectively learn from data while obeying physical laws, resulting in high-level performance for streamflow simulations. From the perspective of process-based modeling, LSTM is a regionalized parameter provider that leverages the autocorrelated nature of its inputs to impose an implicit spatial constraint on the generated parameters.

~~In this study, we used δ HBV1.1p (Song et al., 2024b, 2025) which is an updated version from δ HBV1.0. In this study, we used δ HBV1.1p (Song et al., 2024b, 2025b), which is an updated version of δ HBV1.0 (Feng et al., 2022, 2023b). The main improvement is the addition of a capillary rise module, which enhances the characterization of low flows. Other modifications include the use of three dynamic parameters during the warm-up, training, and test periods, the removal of log-transform normalization for precipitation, and the adoption of NSE as the loss function for model training. The basic equations in δ HBV are as follows:~~

~~. The main improvement is the addition of a capillary rise module, which enhances the characterization of low flows. Other modifications include three modifications to address high-flow simulation challenges: the use of three dynamic parameters (γ , β , k_0) (Song et al., 2025b); the removal of log-transform normalization for precipitation; and the adoption of the~~

normalized squared-error loss function (Frame et al., 2022; Kratzert et al., 2021; Song et al., 2025a, b; Wilbrand et al., 2023). We also maintain dynamic parameters during warm-up periods. Although this provides only marginal benefits and increases computational costs, it yields a more realistic representation and reduces uncertainties associated with initial conditions. The basic equations in δ HBV are as follows:

$$\theta = LSTM_w(\bar{x}, \overline{A_{attr}}) \quad (1)$$

$$Q = HBV(x, \theta) \quad (2)$$

$$W_{opt} = argmin_w(L(Q, Q^*)) \quad (3)$$

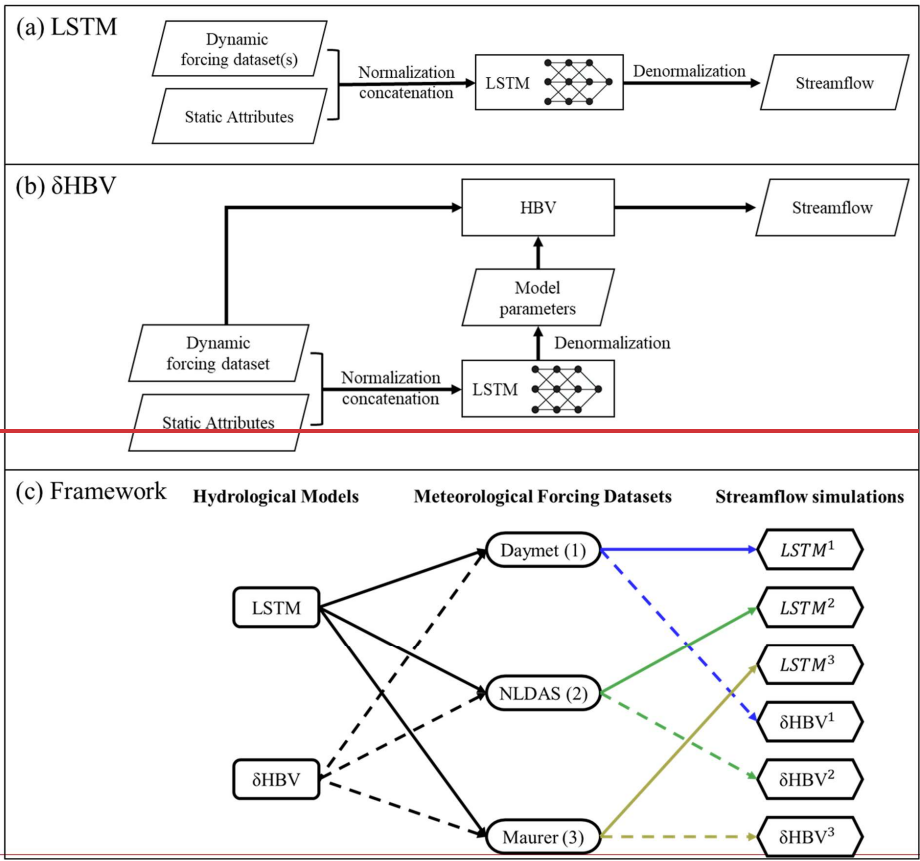
where θ are the dynamic or static physical parameters, w denotes the weights and biases of LSTM, x includes the basin-averaged meteorological forcings, such as precipitation, mean temperature, and potential evapotranspiration, with \bar{x} representing their normalized versions.

Similarly, $\overline{A_{attr}}$ consists of normalized observable basin-averaged attributes, encompassing basin area, topography, climate, soil texture, land cover, and geology (Table C1). Precipitation and mean temperature are from CAMELS, while potential evapotranspiration is calculated based on the Hargreaves (1994) method using mean, maximum, and minimum temperatures along with basin latitudes, all from data described in sect. 2.1. Q and Q^* are the streamflow simulations (model outputs) and observations (as provided in CAMELS), respectively. HBV is implemented on PyTorch so it is programmatically differentiable: all steps store information related to gradient calculations during backpropagation, allowing this model to be trained together with neural networks in an end-to-end fashion. More details about differentiable HBV can be found in previous studies (Feng et al., 2022; Song et al., 2024b). The details of some particularly relevant HBV processes are described in Appendix A.

2.4. Experimental Design

In this study, we trained the two models ~~in highly of very~~ different types (LSTM and δ HBV) using), each with one of three meteorological forcing datasets (Daymet, NLDAS, and Maurer), resulting in six corresponding streamflow simulations (Figure 1c) for each different test scenario (see sect. 2.5 for additional information). The training processes of LSTM and δ HBV followed Kratzert et al. (2021) and Feng et al. (2023b), respectively. Test results and performance metrics for all models are reported for the 531-basin subset that excludes those with areas larger than 2,000 km² or with more than a 10% discrepancy between different basin area calculation methods (Newman et al., 2017).

To generate ensembles, we tested various weighting strategies and ultimately employed averaging to combine the six single-forcing, single-model-type simulations, as it yielded the best performance. To better describe various combinations including cross-model ensembles, these simulations were categorized into six groups (Table 1). A shorthand notation is used throughout the remainder of this work to describe the forcing datasets and ensembles. Daymet, NLDAS, and Maurer are abbreviated as superscripts 1, 2, and 3, respectively. The + symbol is used to group model types being ensembled, while superscript clustering (e.g., ¹² or ¹²³) is used to group the meteorological forcing types being ensembled, with parentheses indicating that the superscripts apply to all model types within. For example, $(LSTM + \delta HBV)^{123}$ could be explicitly written as $LSTM^1 + LSTM^2 + LSTM^3 + \delta HBV^1 + \delta HBV^2 + \delta HBV^3$. To compare two different strategies to utilize the multiple meteorological forcing datasets and to benchmark against the previously highest performance, we additionally trained a single LSTM model using all three forcing datasets as simultaneous inputs as done by Kratzert et al. (2021), -referred to as LSTM^{multi} (the last row in Table 1).



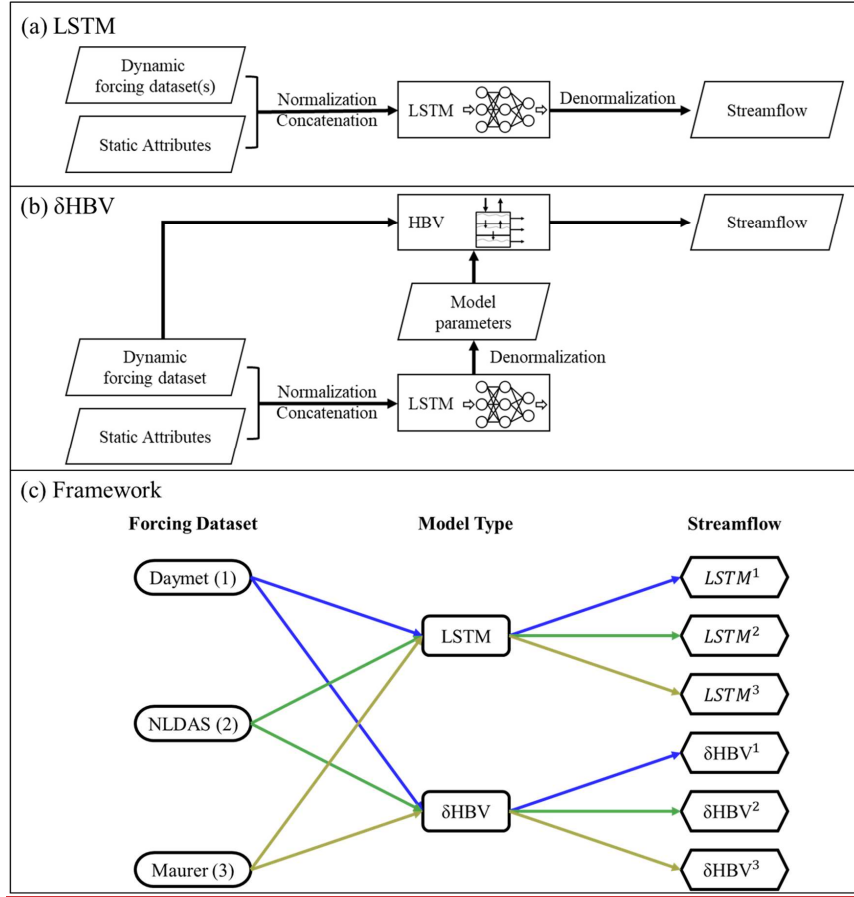


Figure 1. (a) The LSTM structure, (b) the δ HBV structure, and (c) the framework to generate the six individual ensemble members of the streamflow simulations, in which different colors of arrow lines denote the different meteorological forcing datasets (also denoted as 1, 2, 3) and the arrow line styles (solid and dashed lines) indicate the LSTM and δ HBV models, respectively.

Table 1. (a) The six groups of streamflow simulations, and (b) the streamflow simulation via LSTM based on a different strategy ~~to utilize~~, in which three meteorological forcing datasets were combined as a single set of inputs (Kratzert et al., 2021). Superscripts 1, 2, and 3 denote Daymet, NLDAS, and Maurer, respectively. The ensemble across forcings (“ef”) superscript indicates an ensemble of model simulations, each of which uses a different single meteorological forcing, e.g., $LSTM^{12}$ means the average of $LSTM^1$ and $LSTM^2$.

(a) Six Groups of Streamflow Simulations		
Group Name	Group Members	
LSTM	LSTM ¹ , LSTM ² , LSTM ³	
δHBV	δHBV ¹ , δHBV ² , δHBV ³	
LSTM+δHBV	(LSTM+δHBV) ¹ , (LSTM+δHBV) ² , (LSTM+δHBV) ³	
LSTM ^{ef}	LSTM ¹² , LSTM ¹³ , LSTM ²³ , LSTM ¹²³	
δHBV ^{ef}	δHBV ¹² , δHBV ¹³ , δHBV ²³ , δHBV ¹²³ ,	
(LSTM+δHBV) ^{ef}	(LSTM+δHBV) ¹² , (LSTM+δHBV) ¹³ , (LSTM+δHBV) ²³ , (LSTM+δHBV) ¹²³	
(b) Using forcing datasets as simultaneous inputs to an LSTM		
Streamflow Simulation	Model Type	Meteorological Forcing Dataset
LSTM ^{multi}	LSTM	Daymet, NLDAS, Maurer

Formatted: Header

Formatted Table

2.5. Evaluation ~~Scenarios~~scenarios and ~~Criteriacriteria~~

The above cases were comprehensively evaluated for performance in temporal extrapolation (Feng et al., 2022; Kratzert et al., 2018), as well as two types of spatial generalization: prediction in ungauged basins (PUB) (Feng et al., 2023b; Kratzert et al., 2019), and prediction in ungauged regions (PUR) (Feng et al., 2021, 2023b):

- **Temporal Test:** Models were trained using data from all basins and tested across different periods.
- **PUB Test:** Models were trained on randomly selected subsets from all basins and tested on the remaining basins during the same time period.
- **PUR Test:** Different from the PUB test, basins were grouped into continuous regions, one of which was selected to comprise the group of testing basins while the others were used for training.

Formatted: Outline numbered + Level: 1 + Numbering Style: Bullet + Aligned at: 0.33" + Indent at: 0.58"

Temporal generalization is generally considered to be the easiest of these tests. In terms of spatial generalization, which approximates data-sparse scenarios, the PUB test is an example of spatial interpolation, whereas the PUR test involves spatial extrapolation. The PUR test is widely regarded as the most challenging and may therefore produce findings that differ significantly from those in other scenarios. In this study, all basins were divided into 10 spatially stratified groups for the PUB test and 7 fully disjoint regional groups for the PUR test (Table 2) in the same way as Feng et al. (2023b). The spatial extent of the 7 regions for the PUR test is also shown in Figure 3(c1-c2).- Therefore, we conducted 10 rounds for the PUB test and 7 rounds for the PUR test, with a different group held out for testing in each round. Model performance was evaluated after concatenating the test results for all basins.

Formatted: Header

Table 2. Differences of temporal, PUB, and PUR tests.

Test Scenario	Training		Testing	
	Basin	Time	Basin	Time
Temporal	All ^a	1980-1995 ^b	All	1995-2010
PUB	Random nine-tenths	1980-1999	Holdout ^c	1995-1999
PUR	Random six of seven regions	1980-1999	Holdout	1995-1999

Formatted Table

^a δ HBV training followed Feng et al. (2023b) using all 671 CAMELS basins, while LSTM training followed Kratzert et al (2021) using the selected 531-basin subset. Test results and performance metrics for all models are reported for the 531 basins.

^bEach hydrological year spans from October 1st to September 30th of the following year.

^cIn the PUB and PUR tests, models are run for 10 and 7 rounds, respectively, with the group held out for testing changed in each round. The simulation performance was evaluated after concatenating the test results for all basins.

We repeated all the simulations with three different random seeds. Therefore, all the simulations come from a total of $(2 \times 3 + 1) \times (1 + 10 + 7) \times 3$ trained models. The first factor represents the models: two model types (LSTM and δ HBV) trained separately with each of the three forcing datasets, along with $LSTM^{multi}$, a single model instance trained using all three forcing datasets simultaneously. The second factor accounts for the three types of tests (temporal, PUB, and PUR tests), and the last for the three random seeds. With respect to random seeds, we present two variations in the results, which are visually depicted in Figure C1. The results without “seed” as a subscript represent the average metric values from multiple streamflow simulations, each generated from a single model implementation, along with the corresponding uncertainties, visualized using error bars. The results marked with “seed” as a subscript are based on the average of multiple streamflow simulations conducted with different random seeds. In terms of computational cost, training LSTM (30 epochs) and δ HBV (50 epochs) for temporal testing under a single meteorological forcing dataset takes approximately

5 and 21 hours, respectively, using a single NVIDIA Tesla V100 GPU.

We calculated several well-established performance metrics: Nash-Sutcliffe model efficiency coefficient (*NSE*) (Nash and Sutcliffe, 1970), Kling-Gupta model efficiency coefficient (*KGE*) (Kling et al., 2012), percent bias (*PBIAS*), and root-mean-square error (*RMSE*). We also considered *RMSE* values for high (top 2% “peak” flow, *highRMSE*), low (bottom 30% “low” flow, *lowRMSE*), and mid-range (the remaining flow, *midRMSE*) flow conditions (Yilmaz et al., 2008). These metrics were computed for each basin and aggregated into error bars and cumulative density functions (CDFs). ~~Detailed descriptions of these metrics and their calculations are available in Table C2. For brevity, the main text primarily reports NSE values, and other metric values are provided in Appendixes D and E. For brevity, the main text primarily reports NSE values, and other metric values are provided in Appendixes D and E. Furthermore, we use the spread values (Li et al., 2021; Reichle and Koster, 2003) to investigate ensemble variability and explore model complementarity. Detailed descriptions of these metrics and their calculations are available in Table C2.~~

3. Results and ~~Discussion~~discussion

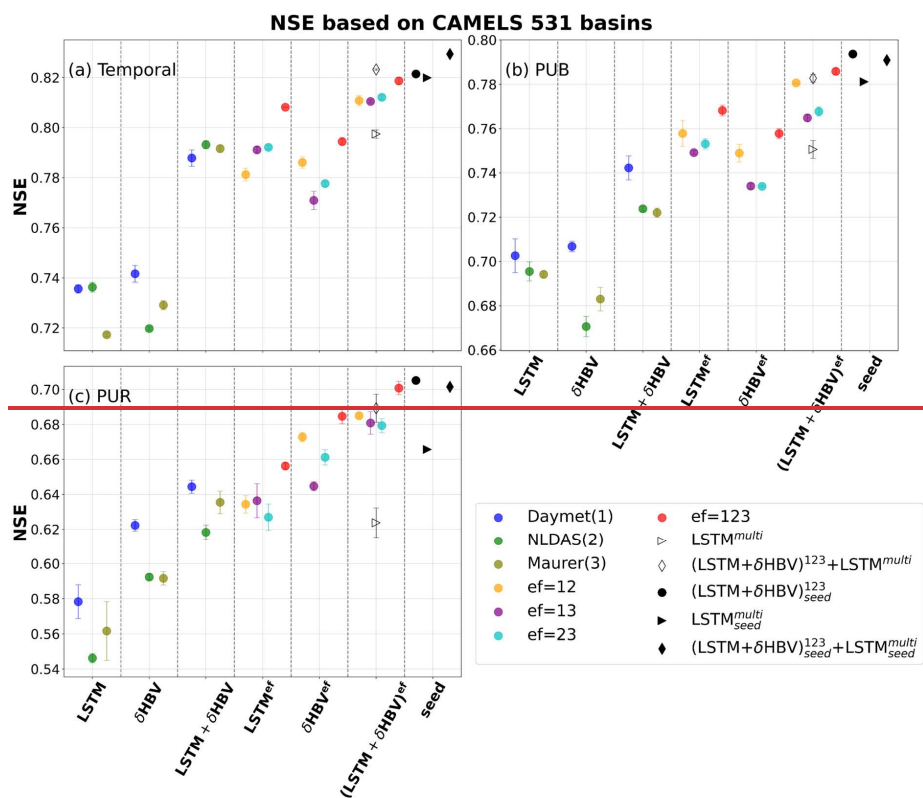
3.1. Temporal extrapolation

For the temporal test, in which models were trained and tested on the same basins but in different time periods, we found that cross-model-type ensembles noticeably surpassed the within-class ensembles when other conditions were the same, with small uncertainties (as shown by the error bars in Figure 2). With a single forcing dataset, the median *NSE* was elevated from ~0.735 for LSTM to ~0.79 with δ HBV added, though δ HBV performance was similar to LSTM (~0.74 under Daymet). Even after LSTM achieved very high performance when its simulations, each derived separately from different meteorological forcing datasets, were ensembled ($ef = 123$, ~0.808), adding δ HBV still improved the results to ~0.818. This

Formatted: Font color: Auto

Formatted: Font color: Auto

finding was robust for all different combinations of the tested meteorological forcing datasets. Conversely, adding LSTM also helped to improve δ HBV ensembles. These results highlight the benefits of the cross-model-type ensemble framework, and indicate distinct simulation features ~~via~~for each model type. LSTM is a data-driven method that has low bias and large variance. ~~Errors with data~~Data errors (Li et al., 2020b), different sampling strategies (Nai et al., 2024), or even different weight initializations (Narkhede et al., 2022) can lead to substantively different outcomes. ~~On the other hand~~Conversely, δ HBV may have a smaller variance but a larger bias due to the fixed HBV formulation (Moges et al., 2016) for some scenarios like low flows (Feng et al., 2023b; Song et al., 2024b) or in basins with significant water uses ~~(Song et al., 2024a)~~(Song et al., 2024a). These errors with varying characteristics from different model classes can partially offset each other in an ensemble. On a side note, δ HBV models seem more reliant on the quality of the forcing data, as shown in Figure 2. δ HBV with the Maurer and NLDAS forcing datasets generally performs worse than it does with Daymet ~~that, which~~ has lower biases. However, even in those cases, the combination of LSTM and δ HBV was still better than LSTM alone, attesting to the robustness of these benefits.



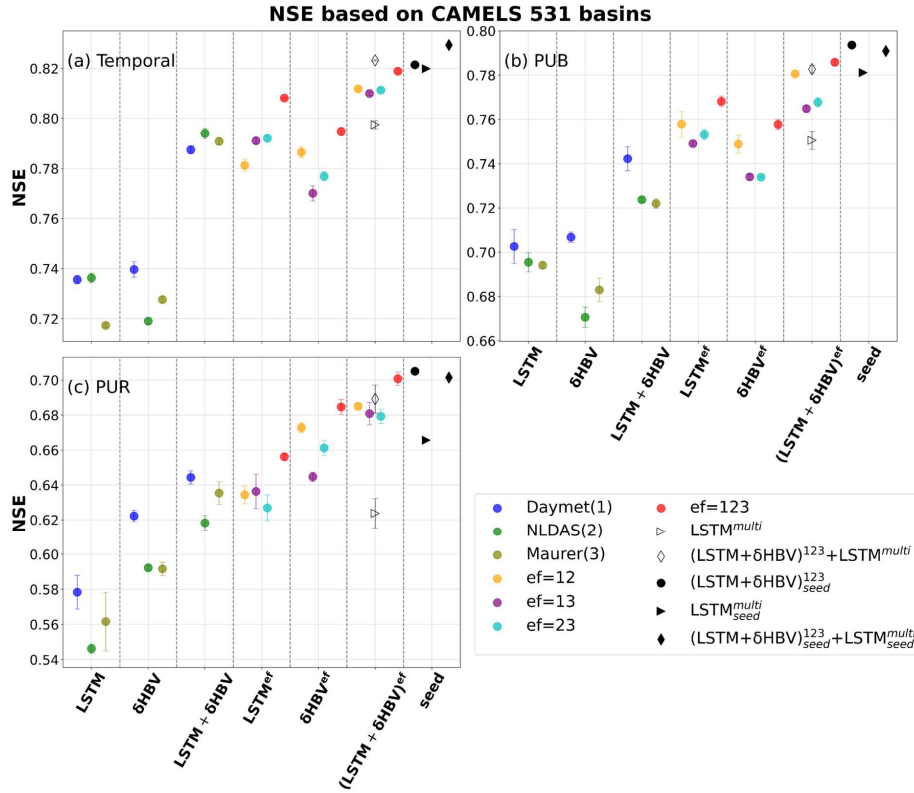


Figure 2. Median NSE values for 531 CAMELS basins, indicating model and ensemble performances for (a) temporal, (b) prediction in ungauged basin (PUB), and (c) prediction in ungauged region (PUR) tests. Different simulations are represented by variously-shaped and -colored points, and are organized by ensemble group, listed along the x-axis: LSTM, δHBV , LSTM+ δHBV , and their “ensemble forcing” counterparts, LSTM^{ef}, δHBV ^{ef}, and (LSTM + δHBV)^{ef}. LSTM^{multi} is a single LSTM model trained directly on all three forcing datasets at once. The superscript “ef” denotes the forcing datasets involved in each ensemble (choices of 1 for Daymet, 2 for NLDAS, and 3 for Maurer), while the “+” connects the model types used within an ensemble. The x-axis group and subscript “seed” indicate that simulation results were averaged based on three different random seeds (see Figure C1). Other points without “seed”, along with their corresponding error bars, are derived from the averages of metrics computed over repeated runs with three different random seeds. The error bar indicates one standard deviation above and below the average value for each simulation.

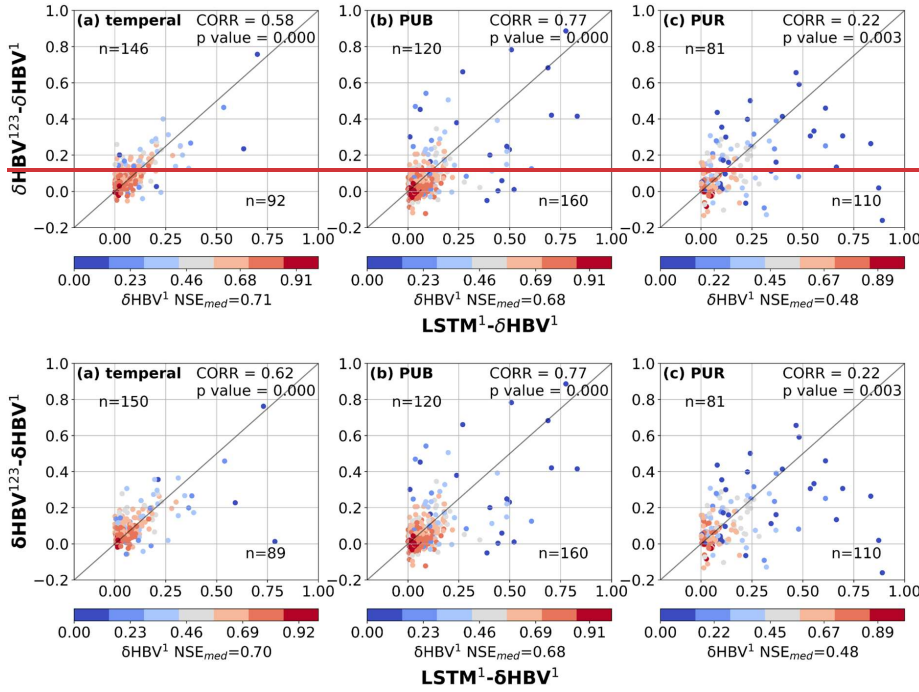


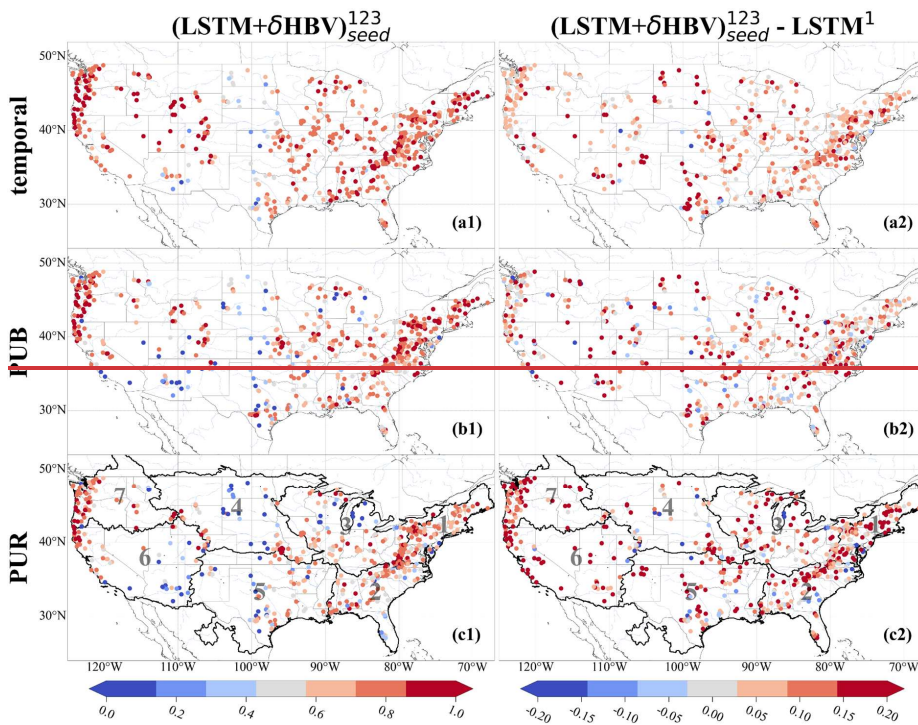
Figure 3. Scatter plots comparing the performance differences between hydrological models for the basins where LSTM outperformed δHBV^1 (the basins where δHBV^1 outperformed are not shown in this plot). The x-axis represents the NSE differences between $LSTM^1$ and δHBV^1 ($LSTM^1 - \delta HBV^1$), while the y-axis shows the NSE differences between δHBV^{123} and δHBV^1 ($\delta HBV^{123} - \delta HBV^1$). Points are color-coded according to the NSE values of δHBV^1 . The correlation coefficient (CORR) and p values between the x-axis values and the y-axis values, along with the median NSE value of δHBV^1 (NSE_{med}) on these basins, are also noted. We note that NSE is not additive and should ~~in-general~~ generally not be subtracted. Here the purpose is only to confirm that basins where LSTM outperforms δHBV^1 also tend to be those that benefit from the ensemble of forcings.

In the lower-performing basins where $LSTM^1$ had advantages over δHBV^1 , the ensemble of meteorological forcings δHBV^{123} also tended to be higher than δHBV^1 (Figure 3), suggesting that forcing quality was a significant reason behind the underperformance of δHBV^1 in these basins. Similar patterns were also observed when analyzing $RMSE_{\delta HBV^2}$ and δHBV^3

values (Figure D1 [and Figure D2](#)). These basins previously contributed to LSTM's cumulative distribution function of NSE diverging from that of δHBV^1 at the low end (Feng et al., 2022). Forcing errors can exist in the form of systematic timing errors, low or high bias for larger events, etc., which can be difficult for the mass-balanced conceptual HBV^1 structure to adapt to these errors. Because the ensemble of forcings tends to suppress the errors in each forcing source, part of the advantages of δHBV^{123} over δHBV^1 can be attributed to reducing forcing bias or timing errors. Since the advantages of LSTM^1 over δHBV^1 also tend to occur with these same basins, this also explains how LSTM^1 surpasses δHBV^1 in some basins with poorer-quality forcings. In contrast to δHBV , LSTM has the innate ability to shift information in time and moderately adjust the input scale. Moving from temporal validation to PUB to PUR scenarios, the advantages of diverse forcing datasets appear to diminish, as evidenced by the decreasing ratio of points above versus below the diagonal line, since the forcing error patterns remembered by LSTM may not generalize well in space (discussed in more detail in sect. 3.2).

Ensembling streamflow simulations from different meteorological forcing datasets demonstrates certain advantages over the previous approach of simultaneously sending multiple forcings into ~~an ML~~ [a data-driven](#) model like LSTM (Kratzert et al., 2021). Ensembling LSTM simulations each using a single forcing dataset (LSTM^{123}) resulted in an NSE value of 0.8082, higher than that of 0.7974 from feeding multiple forcing datasets into a single LSTM (LSTM^{multi}). This difference was more pronounced in the cross-model-type ensemble, after including δHBV , compared to the previous within-class ensemble, and particularly notable for the spatial generalization tests (to be discussed in more detail in Sect. 3.2), with specific metric values provided in Tables D1-D5. These results indicate that the trained LSTM in LSTM^{multi} may be overfitted to the significant redundant information in these three forcing datasets, and that only LSTM cannot fully exploit the information hidden

in the multiple forcing datasets. Training separate ensemble members via different nonlinear hydrological processes, on the other hand, seems to allow different bias features to emerge with separate forcing datasets, accordingly mitigating them during the subsequent ensembling process.



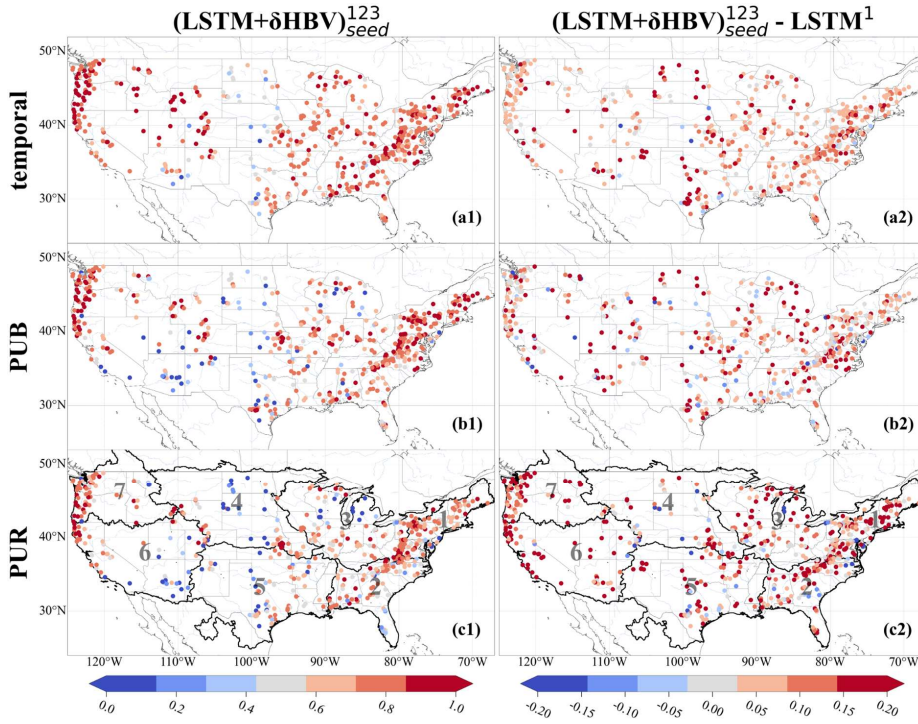


Figure 4. Spatial distributions of NSE values over 531 basins. Subplots are arranged in rows, indicating (a) temporal, (b) PUB, and (c) PUR test results, and columns, denoting (1) NSE values from $(LSTM + \delta HBV)_{seed}^{123}$ and (2) the differences between these NSE values and those of $LSTM^1$ (models using only forcing 1, Daymet). For $LSTM^1$, each NSE value reported was the average of three NSE values from three simulations using three different random seeds. The seven continuous regions used to divide up basins for the PUR test are outlined and numbered in the PUR test maps.

Our most diverse ensemble, $(LSTM + \delta HBV)_{seed}^{123} + LSTM_{seed}^{multi}$, achieved a median NSE value of ~ 0.83 , surpassing the ~ 0.82 benchmark set by $LSTM_{seed}^{multi}$ (Table D4). This advancement was achieved through random seed variation and cross-model-type ensembling. The performance of $(LSTM + \delta HBV)^{123}$ ensemble proved more robust than $LSTM^{multi}$, with only a slight boost when we incorporated random seeds, i.e., $(LSTM + \delta HBV)_{seed}^{123}$. Notably, the derived $(LSTM + \delta HBV)_{seed}^{123}$ ensemble outperformed $LSTM^1$ across almost all basins

(Figure 4). Further incorporation of $LSTM^{multi}$ into this framework, especially when using multiple random seeds, $(LSTM + \delta HBV)_{seed}^{123} + LSTM_{seed}^{multi}$, yielded the best overall performance. Here, the margin over the previous benchmark was small in the temporal test. However, as we will show in sect. 3.2, the previous benchmark, $LSTM_{seed}^{multi}$, lacked robustness, exhibited greater deficiencies in spatial generalization, and negatively impacted ensemble simulations.

When we changed the number of random seeds from 3 to 10, we found that although all model and ensemble performances slightly increased, the gaps between them did not change much (Figure 5 and Table D5). In particular, the gap between $(LSTM + \delta HBV)_{seed}^{123} + LSTM_{seed}^{multi}$ and $(LSTM + \delta HBV)_{seed}^{123}$ or $LSTM_{seed}^{multi}$ remained unchanged. This indicates that the benefits from more random seeds rapidly become marginal, and our results based on 3 random seeds were sufficiently robust. For LSTMs alone, different random seeds displayed higher variation, and ensembling them led to greater improvement than ensembling $(LSTM + \delta HBV)_{seed}^{123}$ with additional random seeds. It was noteworthy that while the $(LSTM + \delta HBV)_{seed}^{123}$ ensemble generally showed the lowest RMSE values, it did not always show the best high flow performance, as indicated by highRMSE (Tables D1-D4). After incorporating the $LSTM_{seed}^{multi}$ variant into $(LSTM + \delta HBV)_{seed}^{123} + LSTM_{seed}^{multi}$, overall RMSE and highRMSE both improved. Nevertheless, this ensemble did not always obtain the best values in other metrics like low flow (lowRMSE) and requires further improvement.

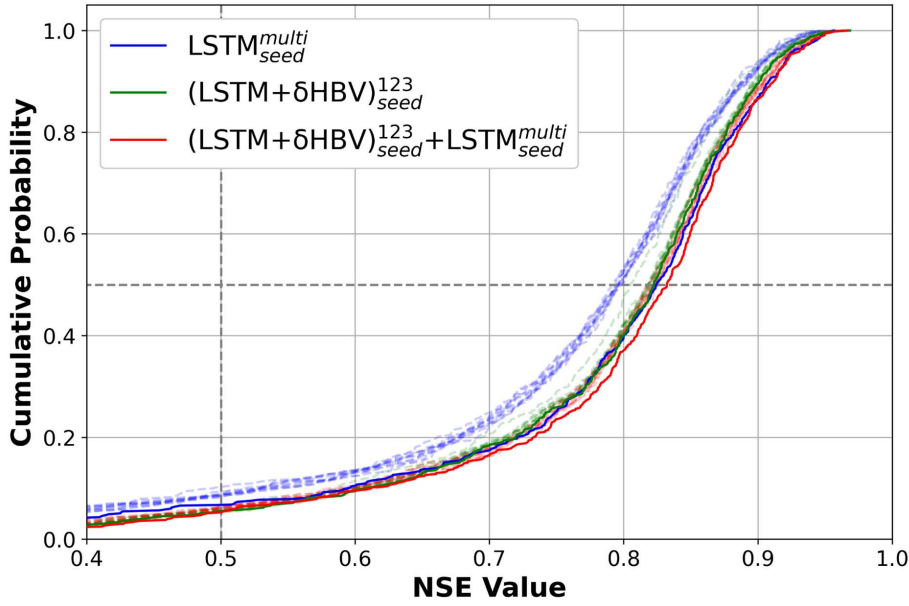
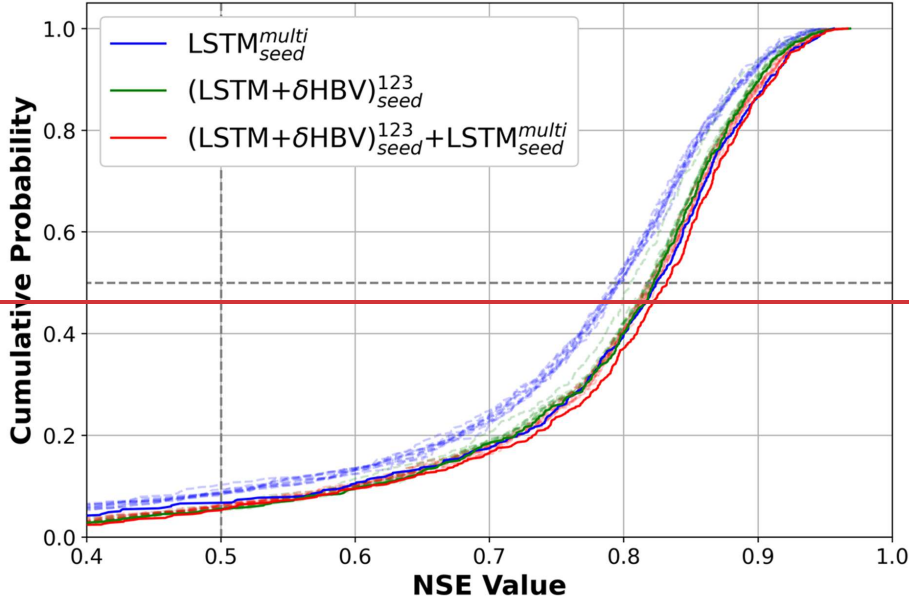


Figure 5. Cumulative distribution function (CDF) curves based on temporal test results for $LSTM^{multi}$, $(LSTM + \delta HBV)^{123}$, and $[(LSTM + \delta HBV)^{123} + LSTM^{multi}]$. The solid lines (with “seed”) denote the results with 10 random seeds while the corresponding dashed and

translucent lines denote the performances of their individual members each based on one random seed.

3.2. Spatial generalization

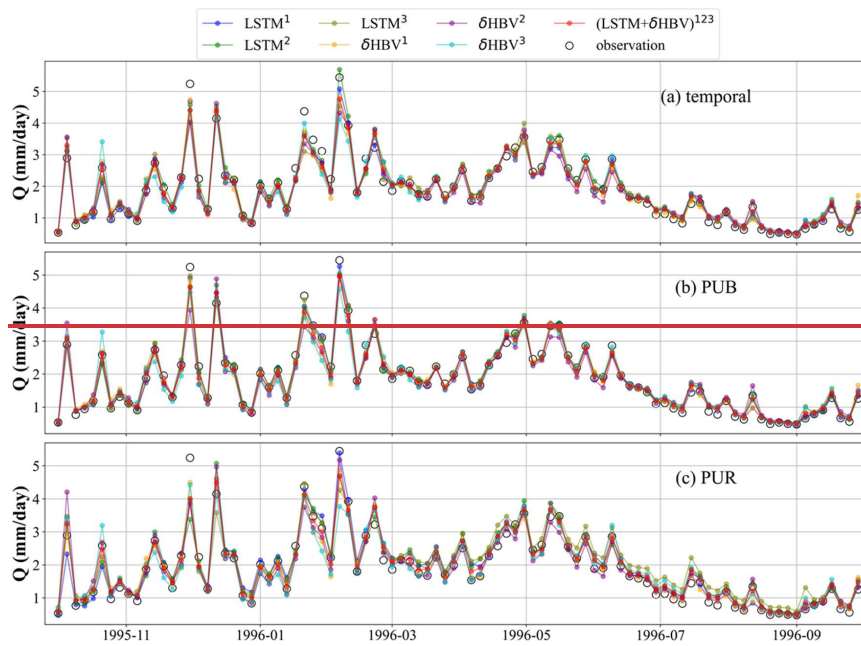
It is clear that cross-model-type ensembling and the incorporation of δ HBV significantly improved prediction in ungauged basins (PUB) or regions (PUR), mitigating the difficulty of spatial generalization (Figure 2b - 2c). In particular, the previous record-holder for temporal test performance, $LSTM_{seed}^{multi}$, incurred large drops in the PUB and PUR tests, once again reminding us of the limitations of LSTM in spatial generalization. Given the same forcings, δ HBV-only individual simulations or ensembles consistently outperformed LSTM-only counterparts in the PUR test. Furthermore, adding δ HBV to the same-model-type LSTM ensembles improved median NSE by 0.02-0.03 for PUB. The role of δ HBV became even more prominent in the harder PUR tests, with an increased gap (0.04-0.07), e.g., $LSTM^{123}$ (median NSE ~ 0.656) and $(LSTM + \delta HBV)^{123}$ (median NSE ~ 0.701). The increased significance of δ HBV is also illustrated by the optimized weights shown in Figure E1. ~~The weights are, which were estimated via~~ using a genetic algorithm ~~using the~~ with streamflow observations ~~during from the~~ test periods. ~~Here the estimated~~ These weights are ~~presented~~ solely ~~used to display~~ illustrate the relative contributions of ~~the~~ different ensemble components. The significantly different spatial distribution patterns of these weights among different test scenarios also indicate the differences among temporal, PUB, and PUR tests (Figures E2-E3). The performance of $(LSTM + \delta HBV)^{123}$ improved compared to $LSTM^{multi}$ whether or not we employed multiple random seeds to form an ensemble. As such, we can conclude that the inclusion of a differentiable process-based model like δ HBV in an ensemble is a systematic way to reduce the risks of failed generalizations of LSTM.

Utilizing a cross-model-type ensemble led to widespread improvements over LSTM-only ensembles, with the exception of a few scattered basins for each temporal (Figure 4-a2), PUB

(Figure 4-b2), and PUR (Figure 4-c2) test. The most significant improvements due to the ensemble were concentrated on the center of the Great Plains along with the midwestern US, while the eastern US was moderately improved, suggesting data uncertainty is a larger issue in the central and midwestern US. The Great Plains have historically had poor performance for all kinds of models (Mai et al., 2022) and even the ensemble model had NSE values of only 0.3-0.4 for many of the basins there, although this still marked significant improvements over LSTM¹ (Figure 4-a2, -b2, -c2). Some western basin NSE values were elevated by more than 0.15 for the temporal test (Figure 4-a2) and even more for PUB and PUR. Meteorological stations are generally sparse on the Great Plains, and an ensemble seems to be an effective way to leverage the different forcing datasets that are available. The poor performances in some basins highlight some remaining deficiencies in current models, which clearly cannot fully consider the heterogeneities of different basins; thus, multiscale formulations that resolve such heterogeneities may have advantages (Song et al., 2024a)(Song et al., 2024a).

To investigate why ensembles outperformed single-model, single-forcing approaches, we compared their temporal, PUB, and PUR test simulation time series against observations for 531 basins (Figure 6). Analysis of averaged hydrological year data revealed that while individual ensemble members using single-source forcing datasets performed similarly for easily simulated periods, they showed significant divergence during challenging periods, particularly peak flows. This divergence stems from distinct systematic errors inherent to different model types and forcing datasets. Notably, LSTM-based simulations alone proved insufficient in generating adequate spread to capture these divergent points. ~~A key finding was that δ HBV exhibited markedly different variation patterns compared to LSTM, and its inclusion substantially increased the ensemble spread.~~ By averaging individual model outputs and stabilizing uncertainties, ensemble simulations achieved effective and robust performance across all conditions, which can be shown via the metric highRMSE and lowRMSE values in

Tables D1-D4. This highlights the critical importance of comprehensive training for each ensemble member, including diverse forcing inputs, full-period model calibration, and rigorous hyperparameter tuning, to enable the development of ensure that each member develops distinct characteristics in their simulation behaviors. These differences allow the ensemble to better represent a range of hydrological responses, particularly under extreme or uncertain conditions. By capturing complementary strengths and compensating for individual weaknesses, such well-trained ensemble members collectively enhance the robustness and accuracy of streamflow simulations, ultimately enhancing ensemble performance.



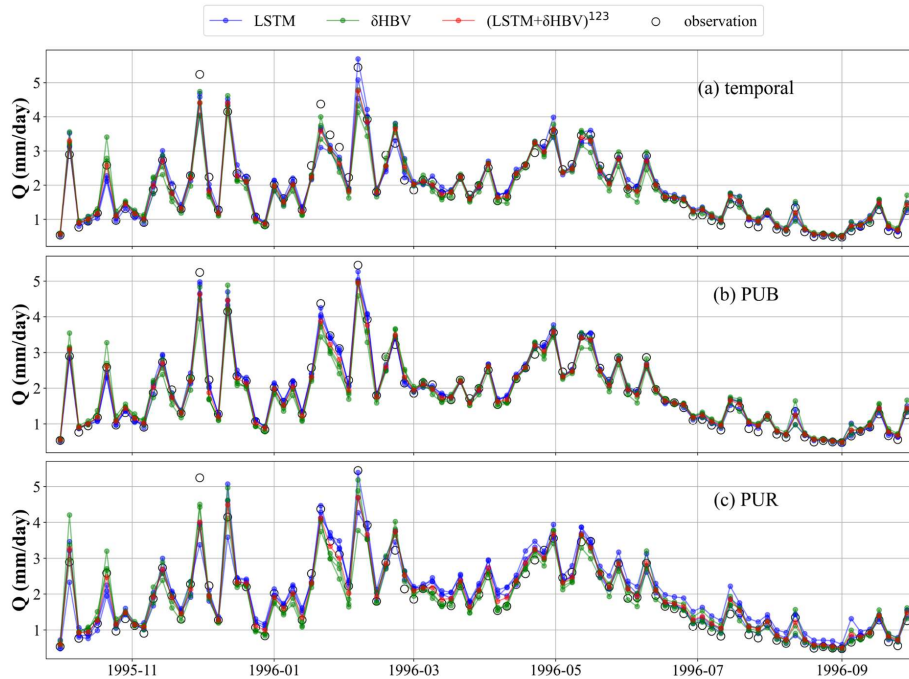


Figure 6. Comparisons between multi-basin-averaged streamflow observations and simulations across 531 basins. The time series points are displayed at four-day intervals for clarity and conciseness. Ensemble members based on the same model (LSTM or δ HBV) but driven by different forcing datasets are shown in the same color to highlight the differences between models more clearly.

Formatted: Font: 12 pt

Formatted: Justified, Indent: First line: 0.33"

3.3 Ensemble variability and robustness analysis

Although δ HBV (median spread 0.61) exhibits lower spreads than LSTM (mean spread 0.72), their combination increases the ensemble spreads, thereby enhancing diversity (Figure 7). This pattern holds across the temporal, PUB, and PUR tests. Ensemble effectiveness depends on the diversity of model behaviors and their distinct error characteristics. Consequently, larger spreads contribute to greater ensemble benefits. Figure D3 further demonstrates that δ HBV+LSTM exhibits larger spreads than LSTM in most basins.

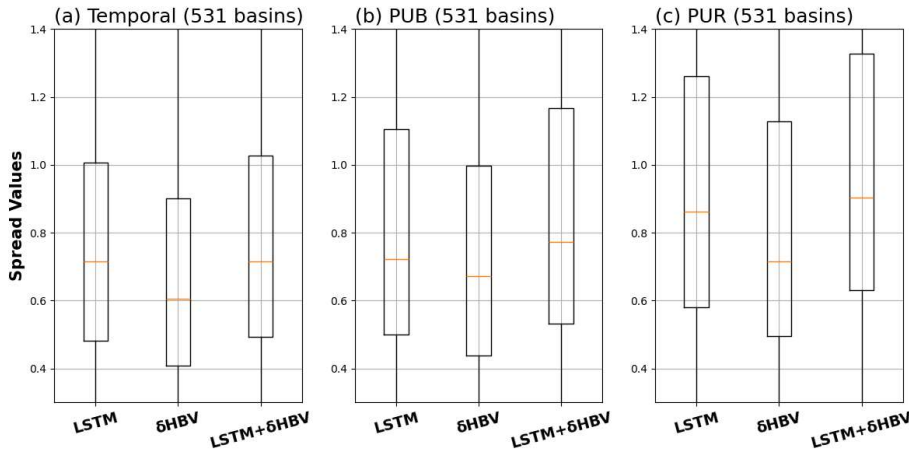


Figure 7. Spread values (Table C2) of each model for LSTM, δ HBV, and LSTM + δ HBV due to different meteorological forcings and random seeds across temporal, PUB, and PUR tests.

As the warming signal is already clear across most basins under any forcing across the periods of simulation (Figure D4), the models' strong performance in the temporal test suggests decent extrapolation capability under warming scenarios. It is often questioned whether data-driven models like LSTM lose accuracy under stronger climate drift, but no substantially warmed dataset is available to test this. Benchmarks suggest LSTM captures 15-year trends well in temporal tests, but less so in data-sparse scenarios (Feng et al., 2023b). Introducing a 10% precipitation perturbation (multiplying precipitation by 1.1) slightly reduced performance for both models as expected (Figures D5a and D5b), but ensemble benefits remained robust across models despite the perturbation.

Training sample size, dynamic parameter choices, and lookback windows all have limited impact on our conclusions. δ HBV shows limited sensitivity to sample size, with similar results when trained on 531 versus 671 basins (Figure D5c). Regarding parameter uncertainties, fixing one δ HBV parameter (k_0) as static increased structural errors and reduced performance (Figure D5d), yet ensemble benefits remained robust. For LSTM, alternative window sizes of 182 and

730 days were tested, with the default 365-day window yielding optimal performance (Figure D5e). Importantly, variations in the lookback window had only minor effects on model performance, underscoring the robustness of ensemble benefits.

3.4. Further discussion

Based on our results, we identified several avenues for future research ~~directions~~. First, while we have explored various weighting strategies and found that averaging yields the best performance yet, we believe that dynamic or adaptive weighting schemes could further enhance performance in future studies. It is also demonstrated by Table E1 that estimated uneven weights can significantly improve simulation performance. Moreover, within specific basins, the estimated weights of different components are often highly imbalanced, as evidenced by the spatial distribution of optimized weights (Figures E2-E3). Some potential feasible ways include using the simulations from these ~~individual-individually~~-trained models as inputs of a data-driven model (Solanki et al., 2025), and making the weight estimation and the ensemble member training simultaneously.

Both LSTM and δ HBV models exhibit limitations in regions with significant anthropogenic impacts like dam presence, as well as arid climatic and significantly heterogeneous geological conditions. These regions are mainly located in the midwestern and western CONUS, where high evaporation conditions (Heidari et al., 2020, Figure 2) and numerous dams (Ryan Bellmore et al., 2017, Figure 1) coincide with complex water use processes (Wada et al., 2016, Figure 11) that current models cannot simulate well. Together, these factors suggest that anthropogenic influence is likely an important driver of poor model performance. Further improvements may include incorporating additional data that capture these factors like capacity-to-runoff ratio (~~Ouyang et al., 2021~~)(Ouyang et al., 2021) or integrating specialized modules, such as reservoirs (Hanazaki et al., 2022; West et al., 2025). Compared with LSTM, δ HBV is more sensitive to precipitation biases. For example, the

differences between δ HBV simulations under different forcing datasets were generally larger than those for LSTM, and δ HBV using the Daymet forcing dataset showed largely better performance than with the other two forcing datasets, which indicates that δ HBV may not be able to fit different forcing datasets well. Therefore, many potential structural optimizations can be implemented to improve δ HBV. Our analysis provided corroborating evidence that forcing error is an important reason why LSTM can outperform δ HBV in the temporal test for some basins, although such patterns may not generalize well in space. A meteorological forcing data correction module can be developed in the future to account for timing and magnitude errors in precipitation. ~~Moreover, ensemble simulations may face challenges when computational resources are limited and calculations are performed sequentially. However, we remain optimistic about these challenges, as the processes can be addressed by leveraging parallel computing with multiple GPUs, benefiting from ongoing advancements in computational power.~~ Ensemble simulations may face challenges when computational resources are constrained, particularly for large-scale or real-time applications. Nevertheless, we remain optimistic about overcoming these challenges due to several promising solutions. These include tailoring the hydrological model by simplifying less relevant components to specific simulation objectives (Clark et al., 2015; Kraft et al., 2022) and cloud-based computing infrastructures that offer scalable, on-demand resource allocation (He et al., 2024; Leube et al., 2013). Importantly, the majority of computational costs are incurred during model training. In practice, ensemble members are typically pre-trained by different research or application groups (Bodnar et al., 2025; Nearing et al., 2024; Song et al., 2025a), enabling direct reuse of these well-trained models and significantly improving computational efficiency.

For this work, we did not create a δ HBV^{multi} model (in the same vein as LSTM^{multi}) using all forcings as an input to a single model, since a similar experiment has already been conducted by Sawadekar et al. (2025). We also did not examine “seed” combinations of a δ HBV^{multi} as

we believed they would not result in a significant performance boost (unlike that seen with LSTM^{multi}), because LSTM has high variability and low bias, while δ HBV has lower variance and potentially higher bias. As a result, random seeds would likely not create large enough perturbations for δ HBV and wouldn't bring the benefits seen with $LSTM_{seed}^{multi}$. To achieve an equivalent perturbation level for δ HBV, it may be necessary to incorporate multiple distinct hydrological models, such as SAC-SMA, PRMS, and GR4J, similar to the approach implemented in the Framework for Understanding Structural Errors (FUSE) (Clark et al., 2008). Work is ongoing to create a combination of a series of differentiable process-based models, which is expected to produce a further improved ensemble with great interpretability. Given the success of cross-model-type ensembles shown in this work, we also encourage further exploration of ensemble simulations involving models with other distinct mechanisms.

4. Summary and ~~Conelusions~~conclusions

This study comprehensively analyzes ensemble combinations of two advanced model types (LSTM and δ HBV), each with distinct mechanisms, for streamflow simulation across 531 basins in the US. Three meteorological forcing datasets (Daymet, NLDAS, and Maurer) are employed to fully capture the characteristics of the two models, ~~and their applications in two different ways are also tested.~~ Their applications are also tested in two distinct ways: (1) by feeding all diverse forcing datasets simultaneously into a single LSTM model, and (2) by ensembling the outputs of multiple LSTM models, each trained separately using a single forcing dataset. The performance of ensemble simulations was evaluated under three distinct testing scenarios (temporal, PUB, and PUR tests), surpassing the previous highest performances. Our findings enhance the understanding of how to effectively utilize diverse model types and multi-source datasets to improve streamflow simulations. The principal conclusions are:

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Header

Formatted: Indent: Hanging: 0.27", Outline numbered + Level: 1 + Numbering Style: 1, 2, 3, ... + Start at: 1 + Alignment: Left + Aligned at: 0.33" + Indent at: 0.6"

- (1) Cross-model-type ensembles (LSTM+ δ HBV) consistently outperformed single-model approaches across all test scenarios, setting new performance benchmarks on the CAMELS dataset. These ensembles demonstrated the complementarity of data-driven (LSTM) and physics-informed (δ HBV) approaches in capturing diverse hydrological behaviors.
- (2) Ensembling models trained on different forcing datasets proved more effective than using multiple forcing datasets as simultaneous inputs to a single model. This suggests that separate training allows each model to capture unique features contained in each forcing dataset, which can then be effectively leveraged in the ensemble.
- (3) δ HBV provided significant benefits to ensemble simulations on spatial generalization. Ensembling LSTM with δ HBV showed increasing benefits as generalization challenges increased, from temporal to spatial interpolation (PUB) to spatial extrapolation (PUR) tests. This underscores the value of physics-informed constraints in improving model transferability to ungauged basins and regions.
- (4) While ensemble methods significantly improved overall performance, they did not fully mitigate consistent deficiencies in certain challenging areas (e.g., regions with high dam density or heterogeneous hydrogeological conditions). This indicates areas for future model development.

These findings have important implications for hydrological modeling and water resources management. The improved accuracy and spatial generalization of our ensemble approach can enhance streamflow predictions, benefiting water resources planning and management, particularly in data-scarce regions. Our results also suggest that future hydrological model development should focus on combining data-driven and physics-based approaches to improve model generalizability across diverse conditions. The superior performance of ensembling models with different forcing datasets over using merged forcings

as a single input highlights the risk of indiscriminately feeding all available data into one data-driven model. While computational demands certainly require consideration, the potential improvements in prediction accuracy offer significant value for both research and operational applications. Future work should focus on refining these ensemble techniques, addressing model limitations in challenging regions, and exploring ensemble implementation in operational settings.

Code and data availability

The source codes and datasets utilized in this study are publicly accessible through the following repositories: The δ HBV modeling framework, including all computational scripts and documentation, is hosted on Zenodo (<https://doi.org/10.5281/zenodo.7091334>) (Feng et al., 2023a), with an updated version and comprehensive software release scheduled upon manuscript acceptance. The implementation of the LSTM architecture is accessible through Zenodo (<https://doi.org/10.5281/zenodo.6326394>) (Kratzert et al., 2022). The CAMELS hydrometeorological dataset, which provides the foundational basin characteristics and time series data used in our analysis, can be obtained via <https://dx.doi.org/10.5065/D6MW2F4D> (Addor et al., 2017; Newman and Clark, 2014).

Author contributions

PL and CS designed the experiments and PL carried them out. YS developed the modified δ HBV code. PL prepared the manuscript with contributions from all co-authors.

Competing interests

Chaopeng Shen and Kathryn Lawson have financial interests in HydroSapient, Inc., a company that could potentially benefit from the results of this research. This interest has been reviewed by the Pennsylvania State University in accordance with its individual conflict of interest policy for the purpose of maintaining the objectivity and the integrity of research. The other authors have no competing interests to declare.

Acknowledgments

PL, CS, and KL were supported by the Office of Biological and Environmental Research of the U.S. Department of Energy (contract no. DESC0016605). PJ and MP were also partially

Formatted: Header

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

supported by California Department of Water Resources Atmospheric River Program Phase III (Grant 4600014294). YS and CS were partially supported by subaward A23-0252-S002 from the Cooperative Institute for Research to Operations in Hydrology (CIROH) through the National Oceanic and Atmospheric Administration (NOAA) Cooperative Agreement (Grant no. NA22NWS4320003).

References

Formatted: Font color: Auto

Aboelyazeed, D., Xu, C., Hoffman, F. M., Liu, J., Jones, A. W., Rackauckas, C., Lawson, K., and Shen, C.: A differentiable, physics informed ecosystem modeling and learning framework for large scale inverse problems: demonstration with photosynthesis simulations, *Biogeosciences*, 20, 2671–2692, <https://doi.org/10.5194/bg-20-2671-2023>, 2023.

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

Aghakouchak, A. and Habib, E.: Application of a Conceptual Hydrologic Model in Teaching Hydrologic Processes, *International Journal of Engineering Education*, 26, 2010.

Bandai, T. and Ghezzehei, T. A.: Physics informed neural networks with monotonicity constraints for Richardson–Richards equation: Estimation of constitutive relationships and soil water flux density from volumetric water content measurements, *Water Resources Research*, 57, e2020WR027642, <https://doi.org/10.1029/2020wr027642>, 2021.

Beek, H. E., van Dijk, A. I. J. M., de Roo, A., Dutra, E., Fink, G., Orth, R., and Schellekens, J.: Global evaluation of runoff from 10 state-of-the-art hydrological models, *Hydrology and Earth System Sciences*, 21, 2881–2903, <https://doi.org/10.5194/hess-21-2881-2017>, 2017.

Beek, H. E., Pan, M., Lin, P., Seibert, J., Dijk, A. I. J. M. van, and Wood, E. F.: Global fully distributed parameter regionalization based on observed streamflow from 4,229

headwater catchments, *Journal of Geophysical Research: Atmospheres*, 125, e2019JD031485, <https://doi.org/10.1029/2019JD031485>, 2020.

Behnke, R., Vavrus, S., Allstadt, A., Albright, T., Thogmartin, W. E., and Radeloff, V. C.: Evaluation of downscaled, gridded climate data for the conterminous United States, *Ecological Applications*, 26, 1338–1351, <https://doi.org/10.1002/15.1061>, 2016.

Bell, V. A. and Moore, R. J.: The sensitivity of catchment runoff models to rainfall data at different spatial scales, *Hydrology and Earth System Sciences*, 4, 653–667, <https://doi.org/10.5194/hess-4-653-2000>, 2000.

Bergström, S.: Development and application of a conceptual runoff model for Scandinavian catchments, PhD Thesis, Swedish Meteorological and Hydrological Institute (SMHI), Norköping, Sweden, 1976.

Bergström, S.: The HBV model – its structure and applications, SMHI, 1992.

Bindas, T., Tsai, W. P., Liu, J., Rahmani, F., Feng, D., Bian, Y., Lawson, K., and Shen, C.: Improving river routing using a differentiable Muskingum-Cunge model and physics-informed machine learning, *Water Resources Research*, 60, e2023WR035337, <https://doi.org/10.1029/2023WR035337>, 2024.

Brunner, M. I., Slater, L., Tallaksen, L. M., and Clark, M.: Challenges in modeling and predicting floods and droughts: A review, *WIREs Water*, 8, e1520, <https://doi.org/10.1002/wat2.1520>, 2021.

Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T., and Hay, L. E.: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between hydrological models, *Water Resources Research*, 44, <https://doi.org/10.1029/2008WR007141>, 2008.

Clark, M. P., Wilby, R. L., Gutmann, E. D., Vano, J. A., Gangopadhyay, S., Wood, A. W., Fowler, H. J., Prudhomme, C., Arnold, J. R., and Brekke, L. D.: Characterizing uncertainty of

the hydrologic impacts of climate change, *Curr Clim Change Rep*, 2, 55–64, <https://doi.org/10.1007/s40641-016-0034-x>, 2016.

Dion, P., Martel, J. L., and Arsenault, R.: Hydrological ensemble forecasting using a multi-model framework, *Journal of Hydrology*, 600, 126537, <https://doi.org/10.1016/j.jhydrol.2021.126537>, 2021.

Feng, D., Fang, K., and Shen, C.: Enhancing streamflow forecast and extracting insights using long-short term memory networks with data integration at continental scales, *Water Resources Research*, 56, e2019WR026793, <https://doi.org/10.1029/2019WR026793>, 2020.

Feng, D., Lawson, K., and Shen, C.: Mitigating prediction error of deep learning streamflow models in large data sparse regions with ensemble modeling and soft data, *Geophysical Research Letters*, 48, e2021GL092999, <https://doi.org/10.1029/2021GL092999>, 2021.

Feng, D., Liu, J., Lawson, K., and Shen, C.: Differentiable, learnable, regionalized process-based models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy, *Water Resources Research*, 58, e2022WR032404, <https://doi.org/10.1029/2022WR032404>, 2022.

Feng, D., Shen, C., Liu, J., Lawson, K., and Beck, H.: differentiable parameter learning (dPL) + HBV hydrologic model, <https://doi.org/10.5281/zenodo.7943626>, 2023a.

Feng, D., Beck, H., Lawson, K., and Shen, C.: The suitability of differentiable, physics-informed machine learning hydrologic models for ungauged regions and climate change impact assessment, *Hydrology and Earth System Sciences*, 27, 2357–2373, <https://doi.org/10.5194/hess-27-2357-2023>, 2023b.

Hanazaki, R., Yamazaki, D., and Yoshimura, K.: Development of a reservoir flood control scheme for global flood models, *JAMES*, 14, e2021MS002944, <https://doi.org/10.1029/2021MS002944>, 2022.

Hargreaves, G. H.: Defining and using reference evapotranspiration, *Journal of Irrigation and Drainage Engineering*, 120, 1132–1139, [https://doi.org/10.1061/\(ASCE\)0733-9437\(1994\)120:6\(1132\)](https://doi.org/10.1061/(ASCE)0733-9437(1994)120:6(1132)), 1994.

Hochreiter, S. and Schmidhuber, J.: Long Short Term Memory, *Neural Computation*, 9, 1735–1780, <https://doi.org/10.1162/neco.1997.9.8.1735>, 1997.

Jiang, S., Zheng, Y., and Solomatine, D.: Improving AI system awareness of geoscience knowledge: Symbiotic integration of physical approaches and deep learning, *Geophysical Research Letters*, 47, e2020GL088229, <https://doi.org/10.1029/2020GL088229>, 2020.

Kling, H., Fuchs, M., and Paulin, M.: Runoff conditions in the upper Danube basin under an ensemble of climate change scenarios, *Journal of Hydrology*, 424–425, 264–277, <https://doi.org/10.1016/j.jhydrol.2012.01.011>, 2012.

Kraft, B., Jung, M., Körner, M., Koirala, S., and Reichstein, M.: Towards hybrid modeling of the global hydrological cycle, *Hydrology and Earth System Sciences*, 26, 1579–1614, <https://doi.org/10.5194/hess-26-1579-2022>, 2022.

Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall Runoff modelling using Long-Short-Term-Memory (LSTM) networks, *Hydrology and Earth System Sciences*, 22, 6005–6022, <https://doi.org/10.17605/OSF.IO/QV5JZ>, 2018.

Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward improved predictions in ungauged basins: Exploiting the power of machine learning, *Water Resources Research*, 55, 11344–11354, <https://doi.org/10.1029/2019WR025448>, 2019.

Kratzert, F., Klotz, D., Hochreiter, S., and Nearing, G. S.: A note on leveraging synergy in multiple meteorological data sets with deep learning for rainfall runoff modeling, *Hydrology and Earth System Sciences*, 25, 2685–2703, <https://doi.org/10.5194/hess-25-2685-2021>, 2021.

Kratzert, F., Gauch, M., Nearing, G., and Klotz, D.: NeuralHydrology — A Python library for Deep Learning research in hydrology, , <https://doi.org/10.5281/zenodo.6326394>, 2022.

Li, P., Zha, Y., Shi, L., Tso, C. H. M., Zhang, Y., and Zeng, W.: Comparison of the use of a physical-based model with data assimilation and machine learning methods for simulating soil water dynamics, *Journal of Hydrology*, 584, 124692, <https://doi.org/10.1016/j.jhydrol.2020.124692>, 2020a.

Li, P., Zha, Y., Tso, C. H. M., Shi, L., Yu, D., Zhang, Y., and Zeng, W.: Data assimilation of uncalibrated soil moisture measurements from frequency domain reflectometry, *Geoderma*, 374, 114432, <https://doi.org/10.1016/j.geoderma.2020.114432>, 2020b.

Li, P., Zha, Y., Zuo, B., and Zhang, Y.: A family of soil water retention models based on sigmoid functions, *Water Resources Research*, 59, e2022WR033160, <https://doi.org/10.1029/2022WR033160>, 2023a.

Li, P., Zha, Y., and Tso, C. H. M.: Reconstructing GRACE derived terrestrial water storage anomalies with in-situ groundwater level measurements and meteorological forcing data, *Journal of Hydrology: Regional Studies*, 50, 101528, <https://doi.org/10.1016/j.ejrh.2023.101528>, 2023b.

Li, P., Zha, Y., Zhang, Y., Michael Tso, C. H., Attinger, S., Samaniego, L., and Peng, J.: Deep learning integrating scale conversion and pedo-transfer function to avoid potential errors in cross-scale transfer, *Water Resources Research*, 60, e2023WR035543, <https://doi.org/10.1029/2023WR035543>, 2024.

Lin, Y., Wang, D., Zhu, J., Sun, W., Shen, C., and Shangguan, W.: Development of objective function-based ensemble model for streamflow forecasts, *Journal of Hydrology*, 632, 130861, <https://doi.org/10.1016/j.jhydrol.2024.130861>, 2024.

Lins, H. F. and Slack, J. R.: Streamflow trends in the United States, *Geophysical Research Letters*, 26, 227–230, <https://doi.org/10.1029/1999GL013086>, 1999.

Liu, J., Rahmani, F., Lawson, K., and Shen, C.: A multiscale deep learning model for soil moisture integrating satellite and in-situ data, *Geophysical Research Letters*, 49,

~~e2021GL096847, <https://doi.org/10.1029/2021GL096847>, 2022.~~

~~Liu, J., Bian, Y., Lawson, K., and Shen, C.: Probing the limit of hydrologic predictability with the Transformer network, Journal of Hydrology, 637, 131389, <https://doi.org/10.1016/j.jhydrol.2024.131389>, 2024.~~

~~Mai, J., Craig, J. R., Tolson, B. A., and Arsenault, R.: The sensitivity of simulated streamflow to individual hydrologic processes across North America, Nat Commun, 13, 455, <https://doi.org/10.1038/s41467-022-28010-7>, 2022.~~

~~Maurer, E. P., Wood, A. W., Adam, J. C., Lettenmaier, D. P., and Nijssen, B.: A long-term hydrologically based dataset of land-surface fluxes and states for the conterminous United States, Journal of Climate, 15, 3237–3251, [https://doi.org/10.1175/1520-0442\(2002\)015<3237:ALTHBD>2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015<3237:ALTHBD>2.0.CO;2), 2002.~~

~~Moges, E., Demissie, Y., and Li, H. Y.: Hierarchical mixture of experts and diagnostic modeling approach to reduce hydrologic model structural uncertainty, Water Resources Research, 52, 2551–2570, <https://doi.org/10.1002/2015WR018266>, 2016.~~

~~Nai, C., Liu, X., Tang, Q., Liu, L., Sun, S., and Gaffney, P. P. J.: A novel strategy for automatic selection of cross-basin data to improve local machine learning-based runoff models, Water Resources Research, 60, e2023WR035051, <https://doi.org/10.1029/2023WR035051>, 2024.~~

~~Narkhede, M. V., Bartakke, P. P., and Sutaone, M. S.: A review on weight initialization strategies for neural networks, Artificial Intelligence Review, 55, 291–322, <https://doi.org/10.1007/s10462-021-10033-z>, 2022.~~

~~Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I — A discussion of principles, Journal of Hydrology, 10, 282–290, [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6), 1970.~~

~~Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz,~~

~~D., Kratzert, F., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev, G., Shenzi, S., Tekalign, T. Y., Weitzner, D., and Matias, Y.: Global prediction of extreme floods in ungauged watersheds, Nature, 627, 559–563, <https://doi.org/10.1038/s41586-024-07145-1>, 2024.~~

~~Newman, A. J. and Clark, M.: A large sample watershed-scale hydrometeorological dataset for the contiguous USA, <https://doi.org/10.5065/D6MW2F4D>, 2014.~~

~~Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., Nearing, G., Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., and Nearing, G.: Benchmarking of a Physically Based Hydrologic Model, Journal of Hydrometeorology, 18, 2215–2225, <https://doi.org/10.1175/JHM-D-18-0114.1>, 2017.~~

~~Newman, A. J., Clark, M. P., Longman, R. J., and Giambelluca, T. W.: Methodological intercomparisons of station-based gridded meteorological products: Utility, limitations, and paths forward, <https://doi.org/10.1175/JHM-D-18-0114.1>, 2019.~~

~~Ouyang, W., Lawson, K., Feng, D., Ye, L., Zhang, C., and Shen, C.: Continental-scale streamflow modeling of basins with reservoirs: Towards a coherent deep learning-based strategy, Journal of Hydrology, 599, 126455, <https://doi.org/10.1016/j.jhydrol.2021.126455>, 2021.~~

~~Paul, P. K., Zhang, Y., Ma, N., Mishra, A., Panigrahy, N., and Singh, R.: Selecting hydrological models for developing countries: Perspective of global, continental, and country scale models over catchment scale models, Journal of Hydrology, 600, 126561, <https://doi.org/10.1016/j.jhydrol.2021.126561>, 2021.~~

~~Rahmani, F., Appling, A., Feng, D., Lawson, K., and Shen, C.: Identifying structural priors in a hybrid differentiable model for stream water temperature modeling, Water Resources Research, 59, e2023WR034420, <https://doi.org/10.1029/2023WR034420>, 2023.~~

~~Sawadekar, K., Song, Y., Pan, M., Beck, H., McCrary, R., Ullrich, P., Lawson, K., and~~

Shen, C.: Hydrology Informed Interpretable Precipitation Data Fusion with a Differentiable Hydrologic Model, <https://doi.org/10.2139/ssrn.4967007>, 25 September 2024.

Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., Baitty Jesi, M., Fenicia, F., Kifer, D., Li, L., Liu, X., Ren, W., Zheng, Y., Harman, C. J., Clark, M., Farthing, M., Feng, D., Kumar, P., Aboelyazeed, D., Rahmani, F., Song, Y., Beck, H. E., Bindas, T., Dwivedi, D., Fang, K., Höge, M., Rackauckas, C., Mohanty, B., Roy, T., Xu, C., and Lawson, K.: Differentiable modelling to unify machine learning and physical models for geosciences, *Nat Rev Earth Environ*, 4, 552–567, <https://doi.org/10.1038/s43017-023-00450-9>, 2023.

Solanki, H., Vegad, U., Kushwaha, A., and Mishra, V.: Improving Streamflow Prediction Using Multiple Hydrological Models and Machine Learning Methods, *Water Resources Research*, 61, e2024WR038192, <https://doi.org/10.1029/2024WR038192>, 2025.

Song, Y., Bindas, T., Shen, C., Ji, H., Knoben, W. J. M., Lonzarich, L., Clark, M. P., Liu, J., van Werkhoven, K., Lemont, S., Denno, M., Pan, M., Yang, Y., Rapp, J., Kumar, M., Rahmani, F., Thébault, C., Sawadekar, K., and Lawson, K.: High-resolution national-scale water modeling is enhanced by multiscale differentiable physics-informed machine learning, <https://doi.org/10.22541/essoar.172736277.74497104/v1>, 26 September 2024a.

Song, Y., Sawadekar, K., Frame, J. M., Pan, M., Clark, M., Knoben, W. J. M., Wood, A. W., Patel, T., and Shen, C.: Improving physics-informed, differentiable hydrologic models for capturing unseen extreme events, <https://doi.org/10.22541/essoar.172304428.82707157/v1>, 2024b.

Song, Y., Knoben, W. J. M., Clark, M. P., Feng, D., Lawson, K., Sawadekar, K., and Shen, C.: When ancient numerical demons meet physics-informed machine learning: adjoint-based gradients for implicit differentiable modeling, *Hydrology and Earth System Sciences*, 28, 3051–3077, <https://doi.org/10.5194/hess-28-3051-2024>, 2024e.

Thornton, P. E., Running, S. W., and White, M. A.: Generating surfaces of daily

meteorological variables over large regions of complex terrain, *Journal of Hydrology*, 190, 214–251, [https://doi.org/10.1016/S0022-1694\(96\)03128-9](https://doi.org/10.1016/S0022-1694(96)03128-9), 1997.

Tsai, W. P., Feng, D., Pan, M., Beck, H., Lawson, K., Yang, Y., Liu, J., and Shen, C.: From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling, *Nat Commun*, 12, 5988, <https://doi.org/10.1038/s41467-021-26107-z>, 2021.

Wang, N., Zhang, D., Chang, H., and Li, H.: Deep learning of subsurface flow via theory-guided neural network, *Journal of Hydrology*, 584, 124700, <https://doi.org/10.1016/j.jhydrol.2020.124700>, 2020.

West, B. D., Maxwell, R. M., and Condon, L. E.: A scalable and modular reservoir implementation for large-scale integrated hydrologic simulations, *Hydrology and Earth System Sciences*, 29, 245–259, <https://doi.org/10.5194/hess-29-245-2025>, 2025.

Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H., Meng, J., Livneh, B., Lettenmaier, D., Koren, V., Duan, Q., Mo, K., Fan, Y., and Moeko, D.: Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of model products, *Journal of Geophysical Research: Atmospheres*, 117, <https://doi.org/10.1029/2011JD016048>, 2012.

Xie, K., Liu, P., Zhang, J., Han, D., Wang, G., and Shen, C.: Physics-guided deep learning for rainfall-runoff modeling by considering extreme events and monotonic relationships, *Journal of Hydrology*, 603, 127043, <https://doi.org/10.1016/j.jhydrol.2021.127043>, 2021.

Yao, L., Libera, D. A., Kheimi, M., Sankarasubramanian, A., and Wang, D.: The roles of climate forcing and its variability on streamflow at daily, monthly, annual, and long-term scales, *Water Resources Research*, 56, e2020WR027111, <https://doi.org/10.1029/2020WR027111>, 2020.

Yilmaz, K. K., Gupta, H. V., and Wagener, T.: A process-based diagnostic approach to model evaluation: Application to the NWS distributed hydrologic model, *Water Resources Research*, 44, <https://doi.org/10/fpvsgb>, 2008.

Yu, D., Yang, J., Shi, L., Zhang, Q., Huang, K., Fang, Y., and Zha, Y.: On the uncertainty of initial condition and initialization approaches in variably saturated flow modeling, *Hydrology and Earth System Sciences*, 23, 2897–2914, <https://doi.org/10.5194/hess-23-2897-2019>, 2019.

Yu, M., Huang, Q., and Li, Z.: Deep learning for spatiotemporal forecasting in Earth system science: a review, *International Journal of Digital Earth*, 17, 2391952, <https://doi.org/10.1080/17538947.2024.2391952>, 2024.

Zhang, Q., Shi, L., Holzman, M., Ye, M., Wang, Y., Carmona, F., and Zha, Y.: A dynamic data-driven method for dealing with model structural error in soil moisture data assimilation, *Advances in Water Resources*, 132, 103407, <https://doi.org/10.1016/j.advwatres.2019.103407>, 2019.

Zounemat Kermani, M., Batelaan, O., Fadaee, M., and Hinkelmann, R.: Ensemble machine-learning paradigms in hydrology: A review, *Journal of Hydrology*, 598, 126266, <https://doi.org/10.1016/j.jhydrol.2021.126266>, 2021.

Appendix A: Detailed processes of HBV employed in this study.

The Hydrologiska Byråns Vattenbalansavdelning (HBV) model (Aghakouchak and Habib, 2010; Beck et al., 2020; Bergström, 1976, 1992) is a simple but effective bucket-type hydrologic model that simulates hydrologic variables including snow water equivalent, soil water, groundwater storage, evapotranspiration, quick flow, baseflow, and total streamflow. In the following texts, we describe these processes in detail by equations, in which uppercase letters indicate state variables, and lowercase letters indicate model parameters. In general, the water balance is developed based on Equation (S1).

$$EP - AE - Q_t = SN + SM + UR + LR + LAKE \quad (S1)$$

where EP is effective precipitation, AE is the actual evapotranspiration, Q_t is the total simulated runoff, SN is snow, SM is soil water storage, UR is the upper reservoir water level, LR is the lower reservoir water level, and $LAKE$ is the lake level (omitted in this study). First, EP is separated into liquid (RN) and solid (SN) components based on the temperature (T) relative to the threshold temperature (tt) as

$$RN = EP \text{ if } T \geq tt \quad (S2)$$

$$SN = EP \text{ if } T < tt \quad (S3)$$

Snow (SN) accumulates in the snowpack (SNP), while the snowmelt (SNM) is calculated using a temperature-dependent melt rate (cfm). The snowmelt (SNM) is limited to the available snowpack (SNP), and any excess melt contributes to meltwater (MW) as

$$SNP = SNP + SN \quad (S4)$$

$$SNM = \begin{cases} SNP & cfm \cdot (T - tt) \geq SNP \\ cfm \cdot (T - tt) & T \geq tt, cfm \cdot (T - tt) < SNP \\ 0 & T < tt \end{cases} \quad (S5)$$

$$MW = MW + SNM \quad (S6)$$

$$SNP = SNP - SNM \quad (S7)$$

Some of this meltwater (MW) refreezes based on a refreezing parameter (cfr) and the temperature difference from the threshold, returning to the snowpack (SNP). The amount of refrozen water is labeled as RFZ .

$$RFZ = \begin{cases} MW & cfr \cdot cfm \cdot (tt - T) \geq MW \\ cfr \cdot cfm \cdot (tt - T) & T < tt, cfr \cdot cfm \cdot (tt - T) < MW \\ 0 & T \geq tt \end{cases} \quad (S8)$$

$$SNP = SNP + RFZ \quad (S9)$$

$$MW = MW - RFZ \quad (S10)$$

The remaining meltwater (MW) that exceeds the snowpack's holding capacity (cwh) contributes

Formatted: Header

Formatted: Font color: Auto

Formatted Table

Formatted Table

Formatted Table

Formatted Table

Formatted: Header

922 to soil infiltration (IF), and the rest remains in the meltwater (MW) storage as

$$IF = \begin{cases} MW - cwh * SNP & MW - cwh * SNP \geq 0 \\ 0 & MW - cwh * SNP < 0 \end{cases} \quad (S11)$$

Formatted Table

$$MW = MW - IF \quad (S12)$$

923 The fraction of soil moisture relative to the field capacity (fc) determines the soil wetness,
924 which modulates the amount of water recharged into the soil (SP). Then soil moisture (SM) is
925 updated based on the infiltration of meltwater (IF), rain (RN), and the amount of recharged
926 water (SP) as

$$SP = \left(\frac{SM}{fc} \right)^\beta \cdot (IF + RN) \quad (S13)$$

Formatted Table

$$SM = SM + IF + RN - SP \quad (S14)$$

927 The excess water, above the field capacity (IF_{dir}), is calculated and subsequently removed from
928 the soil moisture storage as

$$IF_{dir} = \begin{cases} SM - fc & \text{if } SM \geq fc \\ 0 & \text{if } SM < fc \end{cases} \quad (S15)$$

Formatted Table

$$SM = SM - IF_{dir} \quad (S16)$$

929 Actual evapotranspiration (AE) is determined by an evaporation factor (PEC), which depends
930 on the soil moisture, a shape parameter (λ), a parameter (lp), and field capacity (fc) for
931 evapotranspiration. This factor limits the actual evapotranspiration (AE) to both the potential
932 evapotranspiration (PE) and the available soil moisture.

$$PEC = \begin{cases} \left(\frac{SM}{lp \cdot fc} \right)^\lambda & \text{if } 0 \leq \left(\frac{SM}{lp \cdot fc} \right)^\lambda < 1 \\ 0 & \text{if } S \left(\frac{SM}{lp \cdot fc} \right)^\lambda < 0 \\ 1 & \text{if } S \left(\frac{SM}{lp \cdot fc} \right)^\lambda \geq 1 \end{cases} \quad (S17)$$

Formatted Table

$$AE = \begin{cases} PE \cdot PEC & \text{if } SM \geq PE \cdot PEC \\ SM & \text{if } SM < PE \cdot PEC \end{cases} \quad (S18)$$

$$SM = SM - AE \quad (S19)$$

933 Capillary rise (CP) from the lower soil zone (SLZ) is governed by a parameter (c), which
934 determines the amount of water moving upward based on the soil moisture content. This
935 capillary flow replenishes the soil moisture, while groundwater interactions occur through
936 recharge processes in the upper (SUZ) and lower (SLZ) groundwater zones.

$$CP = \begin{cases} SLZ & \text{if } SLZ < c \cdot SLZ \cdot \left(1 - \frac{SM}{fc}\right) \\ c \cdot SLZ \cdot \left(1 - \frac{SM}{fc}\right) & \text{if } SLZ \geq c \cdot SLZ \cdot \left(1 - \frac{SM}{fc}\right) \end{cases} \quad (S20)$$

$$SM = SM + CP \quad (S21)$$

$$SLZ = \begin{cases} SLZ - CP & \text{if } SLZ \geq CP \\ 0 & \text{if } SLZ < CP \end{cases} \quad (S22)$$

937 Excess recharge (SP and IF_{dir}) from the soil enters the upper zone, where it either percolates
 938 to the lower zone ($PERC$) based on a constant rate (prc) or contributes to direct runoff (Q_0)
 939 when it exceeds the upper zone threshold (uzl). The generated flow is modeled using
 940 parameters (k_0, k_1, k_2) governing flow from the upper and lower zones. Each of these flows
 941 contributes to runoff (Q_0, Q_1, Q_2), and their respective contributions to streamflow (Q_t) are
 942 modeled over time.

$$SUZ = SUZ + SP + IF_{dir} \quad (S23)$$

$$PERC = \begin{cases} prc & \text{if } SUZ \geq prc \\ SUZ & \text{if } SUZ < prc \end{cases} \quad (S24)$$

$$SUZ = SUZ - PERC \quad (S25)$$

$$Q_0 = \begin{cases} k_0 \cdot (SUZ - uzl) & \text{if } SUZ \geq uzl \\ 0 & \text{if } SUZ < uzl \end{cases} \quad (S26)$$

$$SUZ = SUZ - Q_0 \quad (S27)$$

$$Q_1 = SUZ \cdot k_1 \quad (S28)$$

$$SUZ = SUZ - Q_1 \quad (S29)$$

$$SLZ = SLZ + PERC \quad (S30)$$

$$Q_2 = SLZ \cdot k_2 \quad (S31)$$

$$SLZ = SLZ - Q_2 \quad (S32)$$

$$Q_t = Q_0 + Q_1 + Q_2 \quad (S33)$$

943
 944 Finally, a routing module (Feng et al., 2022) is used to process Q_t to produce the final
 945 streamflow output (Q_t^*). This module with two parameters ($\theta_\alpha, \theta_\tau$) assumes a gamma function
 946 for the unit hydrograph and convolves the unit hydrograph with the runoff as,
 947

$$Q_t^* = \int_0^{tmax} \xi(s; \theta_\alpha, \theta_\tau) \cdot Q(t - s) ds \quad (S34)$$

Formatted: Header

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted Table

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted Table

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted Table

Formatted Table

|

|

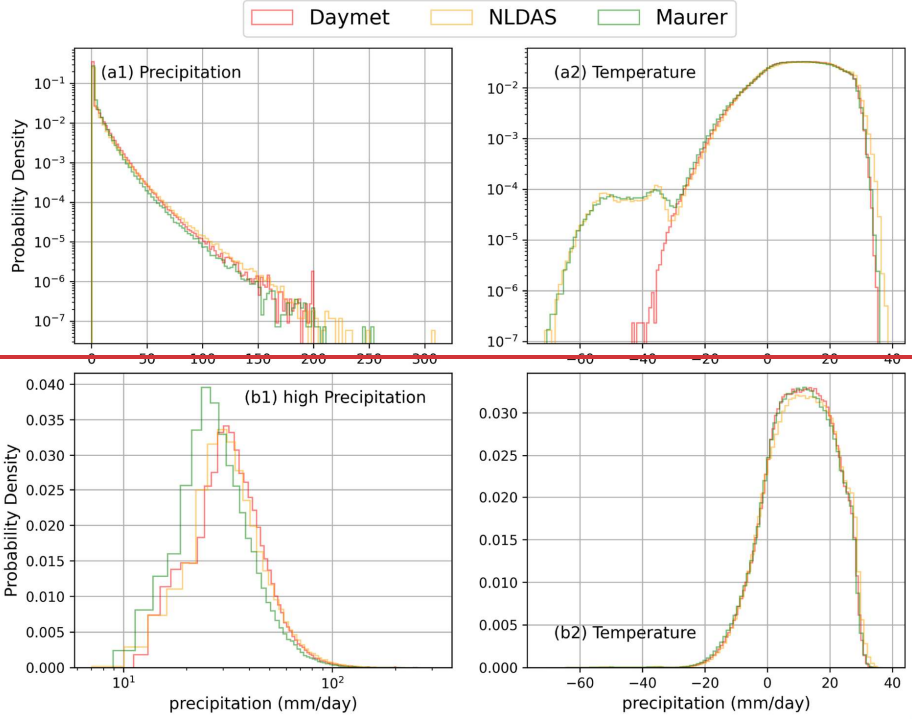
948

949

950

$$\xi(s; \theta_\alpha, \theta_\tau) = \frac{1}{\Gamma(\theta_\alpha)\theta_\tau^{\theta_\alpha}} t^{\theta_\alpha-1} e^{-\frac{t}{\theta_\tau}} \tag{S35}$$

951 **Appendix B: Illustrated differences among the three meteorological forcing datasets**



952

Formatted: Header

Formatted: Font color: Auto

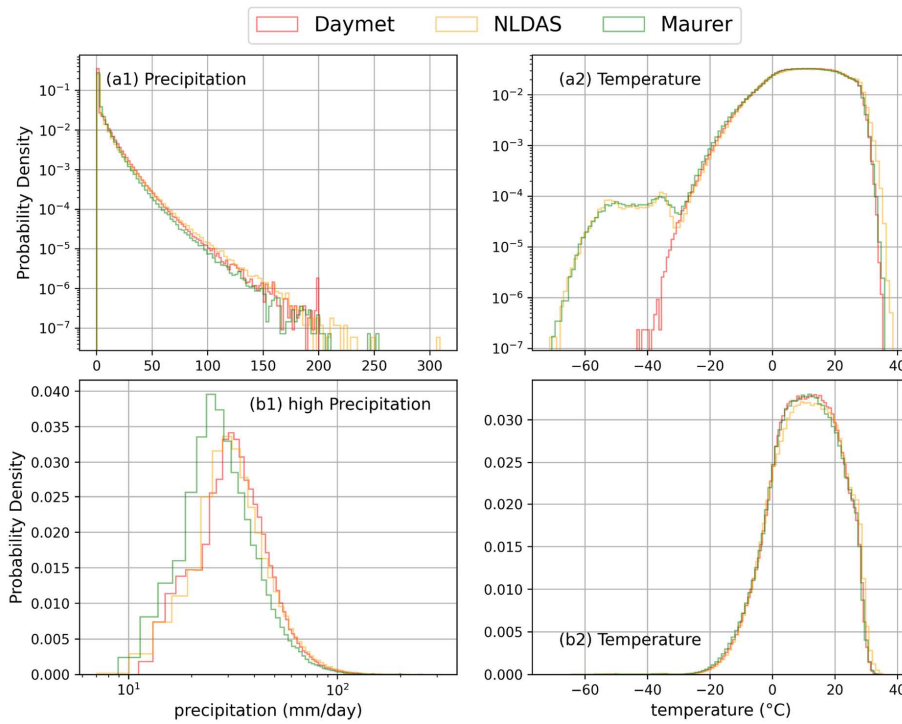


Figure B1. Probability density distributions of precipitation and temperature across three meteorological forcing datasets.

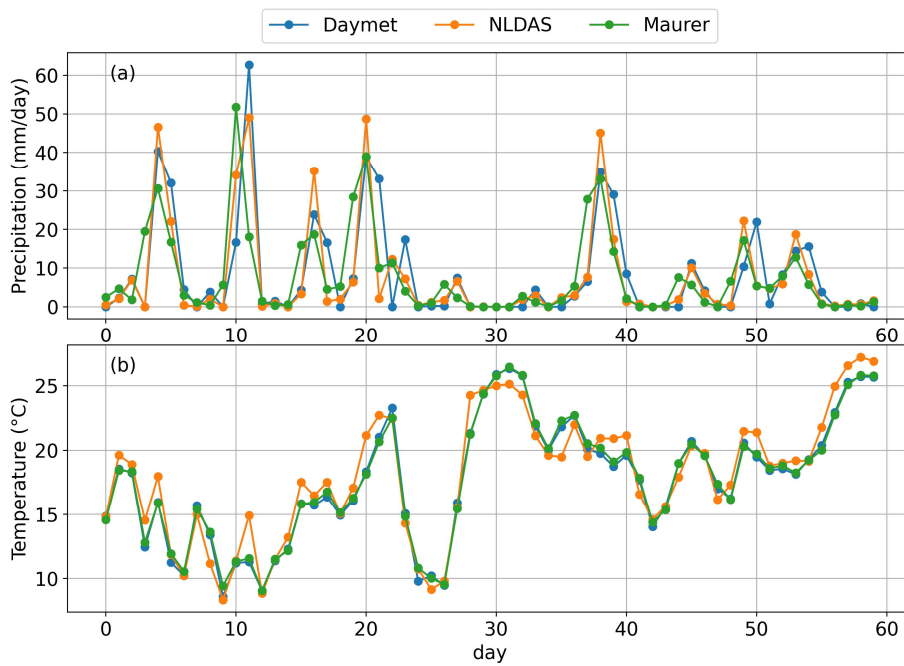
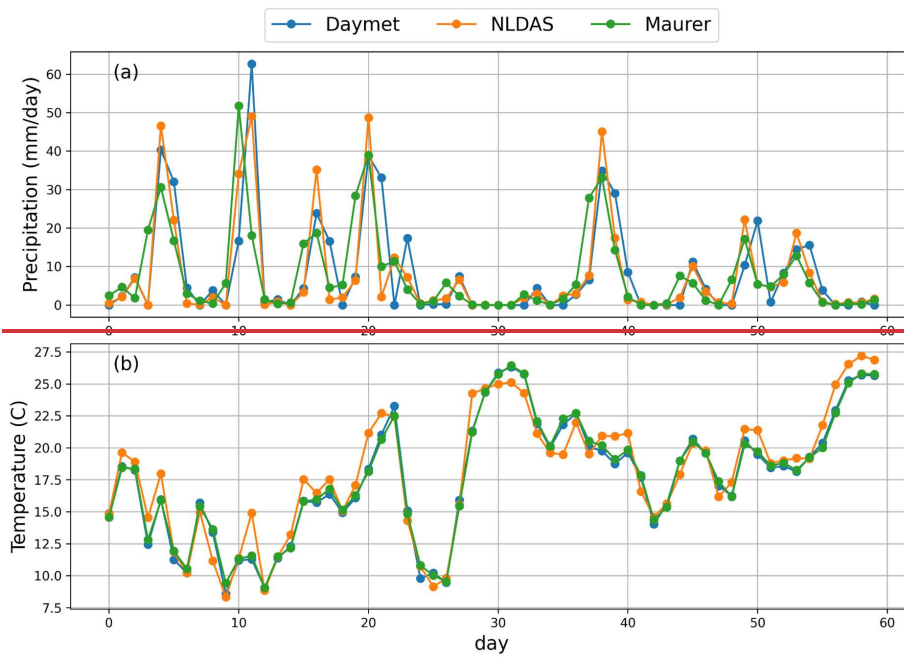


Figure B2. Illustrated temporal variations of precipitation and temperature in a basin across

|

961 *three meteorological forcing datasets.*

962

963

Formatted: Header

Appendix C: Details of model inputs, ensemble frameworks, and evaluations

Table C1. Full names for the abbreviations of dynamic data (all but streamflow are “forcings”) and static basin attributes used as model inputs and outputs. All variables and their values are provided in the CAMELS dataset (Addor et al., 2017) except for the NLDAS and Maurer daily temperature extrema, which are from Kratzert et al. (2021). Potential evapotranspiration and normalized streamflow were calculated ~~for the purposes of in~~ this work, using CAMELS data. The number in parentheses specifies model usage: 1 denotes use in the LSTM model, and 2 denotes use in the δ HBV model.

Type	Abbreviation	Full name	Unit
Dynamic data	prcp <u>(1,2)</u>	Precipitation	mm/day
	pet <u>(2)</u>	Potential evapotranspiration (calculated in this work using the Hargreaves equation and CAMELS data)	mm/day
	tmean <u>(2)</u>	Mean air temperature	°C
	tmax <u>(1)</u>	Maximum air temperature	°C
	tmin <u>(1)</u>	Minimum air temperature	°C
	dayl	Day length	s/day
	srad <u>(1)</u>	Shortwave radiation	W/m ²
	vp <u>(1)</u>	Water vapor pressure	pa
	q_vol	Volumetric streamflow	ft ³ /s
	q <u>(1,2)</u>	Streamflow normalized by basin area (q_vol / area_gages2)	mm/day
Static basin attributes	p_mean <u>(1,2)</u>	Mean daily precipitation	mm/day
	pet_mean <u>(1,2)</u>	Mean daily potential evapotranspiration	mm/day
	p_seasonality <u>(2)</u>	Seasonality and timing of precipitation	-
	frac_snow <u>(1,2)</u>	Fraction of precipitation falling as snow	-
	aridity <u>(1,2)</u>	Rate of mean values of potential evapotranspiration <u>and</u> precipitation	-
	high_prec_freq <u>(1,2)</u>	Frequency of high precipitation days	days/year
	high_prec_dur <u>(1,2)</u>	Average duration of high precipitation events	days
	low_prec_freq <u>(1,2)</u>	Frequency of dry days	days/year

low_prec_dur (1,2)	Average duration of dry periods	days
elev_mean (1,2)	Catchment mean elevation	m
slope_mean (1,2)	Catchment mean slope	m/km
area_gages2 (1,2)	Catchment area (GAGES-II estimate)	km ²
frac_forest (1,2)	Fraction of catchment area having land cover identified as forest	-
lai_max (1,2)	Maximum monthly mean of the leaf area index	-
lai_diff (1,2)	Difference between the maximum and minimum monthly mean of the leaf area index	-
gvf_max (1,2)	Maximum monthly mean of the green vegetation	-
gvf_diff (1,2)	Difference between the maximum and minimum monthly mean of the green vegetation fraction	-
dom_land_cover_frac (2)	Fraction of the catchment area associated with the dominant land cover	-
dom_land_cover (2)	Dominant land cover type	-
root_depth_50 (2)	Root depth at 50 th percentile, extracted from a root depth distribution based on the International Geosphere-Biosphere Programme (IGBP) land cover	m
soil_depth_pelletier (1,2)	Depth to bedrock	m
soil_depth_statsgso (1,2)	Soil depth	m
soil_porosity (1,2)	Volumetric soil porosity	-
soil_conductivity (1,2)	Saturated hydraulic conductivity	cm/hr
max_water_content (1,2)	Maximum water content	m

sand_frac (1,2)	Fraction of soil which is sand	-
silt_frac (1,2)	Fraction of soil which is silt	-
clay_frac (1,2)	Fraction of soil which is clay	-
geol_class_1st (2)	Most common geologic class in the catchment basin	-
geol_class_1st_frac (2)	Fraction of the catchment area associated with its most common geologic class	-
geol_class_2nd (2)	Second most common geologic class in the catchment basin	-
geol_class_2nd_frac (2)	Fraction of the catchment area associated with its 2nd most common geologic class	-
carbonate_rocks_frac (1,2)	Fraction of the catchment area as carbonate sedimentary rocks	-
geol_porosity (2)	Subsurface porosity	-
geol_permeability (1,2)	Subsurface permeability	m ²

972

973

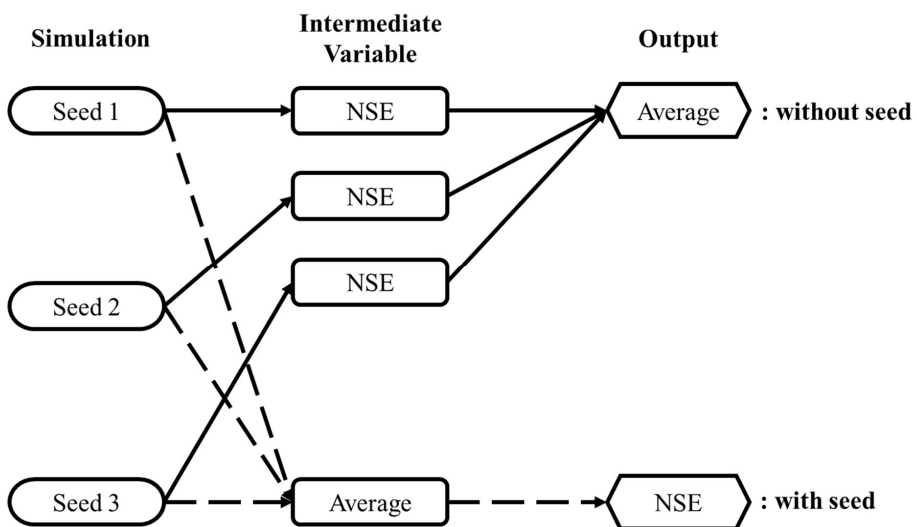
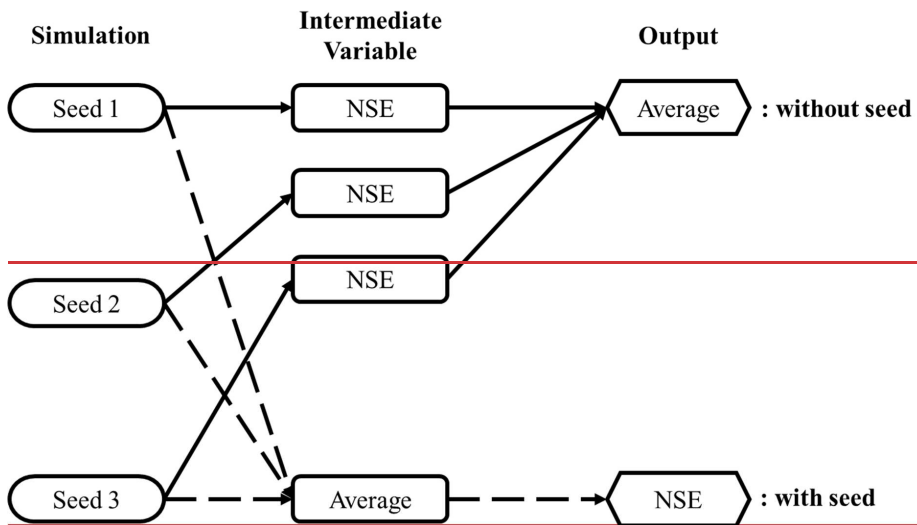


Figure C1. Ensemble frameworks to generate metrics for ensembles named without (solid arrows) and with (dashed arrows) "seed" as a subscript.

Table C2. Evaluation metrics.

Statistic	Equation*	Range	Optimal Value
NSE	$NSE = 1 - \frac{\sum_{i=1}^n (O_i - S_i)^2}{\sum_{i=1}^n (O_i - \mu_o)^2}$	$-\infty$ to 1.0	1.0
KGE	$KGE = 1 - \frac{\sqrt{(r-1)^2 + (\beta-1)^2 + (\gamma-1)^2}}{\beta = \frac{\mu_s}{\mu_o}, \gamma = \frac{CV_s}{CV_o} = \frac{\sigma_s/\mu_s}{\sigma_o/\mu_o}}$	$-\infty$ to 1.0	1.0
PBIAS	$\frac{\sum_{i=1}^n (O_i - S_i)}{\sum_{i=1}^n O_i} \times 100$	$-\infty$ to ∞	0.0
RMSE	$\sqrt{\frac{1}{n} \sum_{i=1}^n (O_i - S_i)^2}$	0.0 to ∞	0.0
spread	$\sqrt{\frac{1}{n} \frac{1}{e} \sum_{i=1}^n \sum_{j=1}^e (S_{i,j} - \mu_{s,i})^2}$	0.0 to ∞	None

* S is a streamflow simulation; O is the corresponding observation; n is the number of total S or O ; e is the number of ensemble members; r is the linear Pearson correlation between S and O ; β is the mean bias; and γ is the variability bias. The mean and standard deviation of simulations are denoted as μ_s and σ_s , respectively, and μ_o and σ_o are the mean and standard deviation of the observations.

Formatted: Header

Inserted Cells

Inserted Cells

Table C2 (continued). Evaluation metrics.

Statistic	Range	Optimal Value
NSE	$-\infty$ to 1.0	1.0
KGE	$-\infty$ to 1.0	1.0
PBIAS	$-\infty$ to ∞	0.0
RMSE	0.0 to ∞	0.0

* S is a streamflow simulation; O is the corresponding observation; n is the number of total S or O ; r is the linear Pearson correlation between S and O ; β is the mean bias; and γ is the variability bias. The mean and standard deviation of simulations are denoted as μ_S and σ_S , respectively, and μ_O and σ_O are the mean and standard deviation of the observations.

|

|990

Formatted: Header

Appendix D: Additional details on model performance

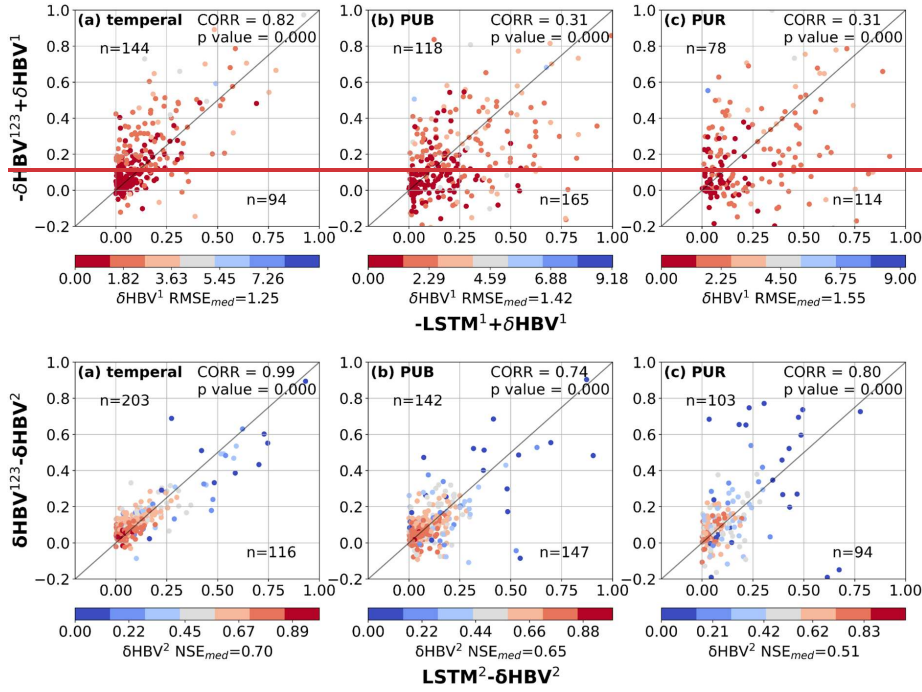


Figure D1. Scatter plots comparing the performance differences between hydrological models for the basins where LSTM outperformed δHBV (the basins where δHBV outperformed are not shown in this plot). The x-axis represents the $RMSENSE$ differences between $LSTM^1$ and δHBV^1 ($LSTM^1 - \delta HBV^1$), while the y-axis shows the $RMSENSE$ differences between δHBV^{123} and δHBV^1 ($\delta HBV^{123} - \delta HBV^1$). Points are color-coded according to the $RMSENSE$ values of δHBV^1 . The correlation coefficient (CORR) and p values between the x-axis values and the y-axis values, along with the median NSE values of δHBV^1 ($RMSENSE_{med}$) on these basins, are also noted.

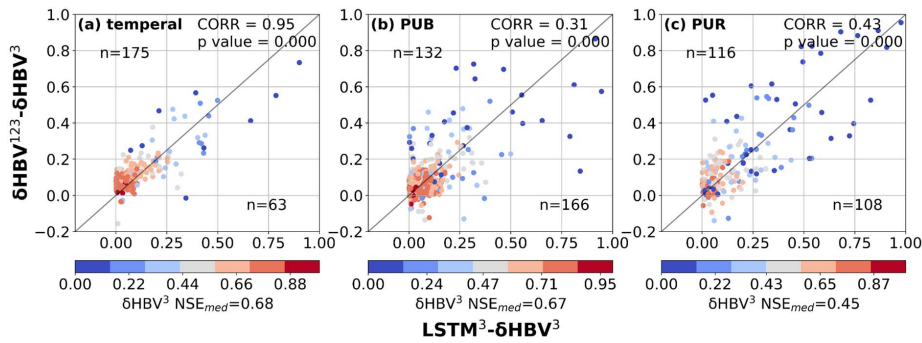


Figure D2. Scatter plots comparing the performance differences between hydrological models for the basins where LSTM outperformed δHBV (the basins where δHBV outperformed are not shown in this plot). The x-axis represents the NSE differences between $LSTM^3$ and δHBV^3 ($LSTM^3 - \delta HBV^3$), while the y-axis shows the NSE differences between δHBV^{123} and δHBV^3 ($\delta HBV^{123} - \delta HBV^3$). Points are color-coded according to the NSE values of δHBV^3 . The correlation coefficient (CORR) and p values between the x-axis values and the y-axis values, along with the median NSE value of δHBV^3 (NSE_{med}) on these basins, are also noted.

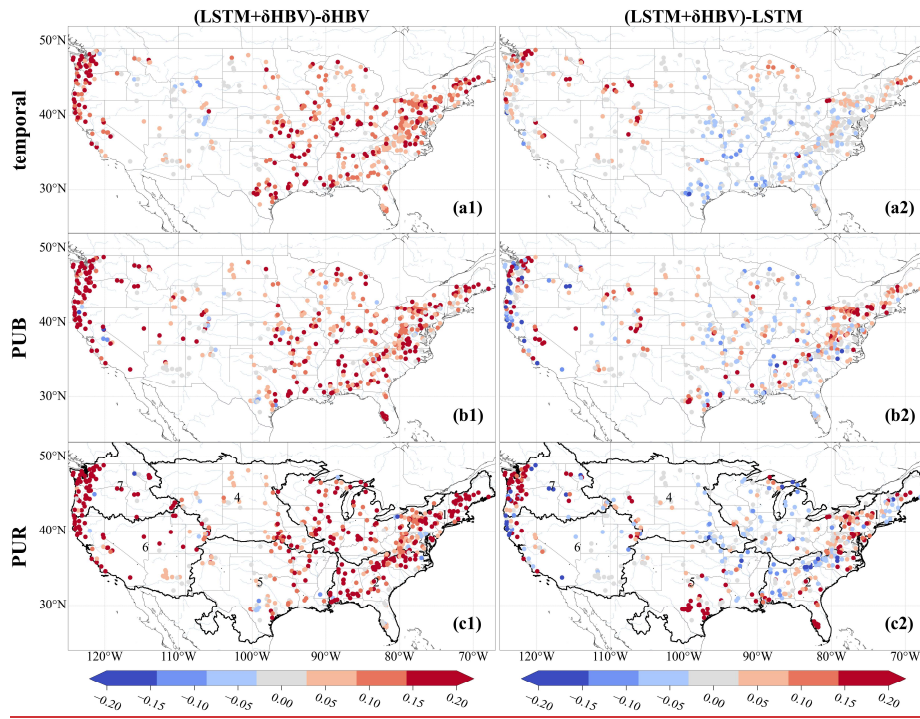


Figure D3. Spatial distributions of model spread values increase from δHBV and LSTM to the LSTM+ δHBV ensemble across temporal, PUB, and PUR tests.

Relative Temperature Differences: (Test – Training) / Training

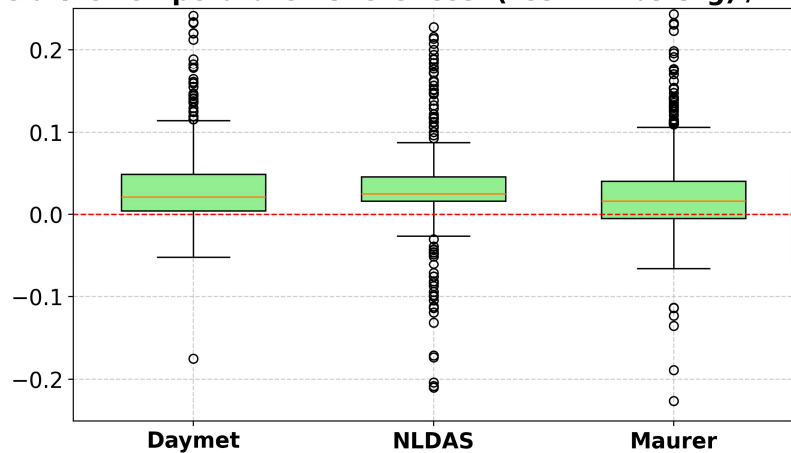


Figure D4. Boxplot of relative temperature differences between the test and training periods, calculated as (Test – Training) / Training. Each box represents the distribution of normalized temperature changes across basins for a specific meteorological forcing dataset: Daymet, NLDAS, and Maurer. Positive values indicate warming in the test period relative to the training period.

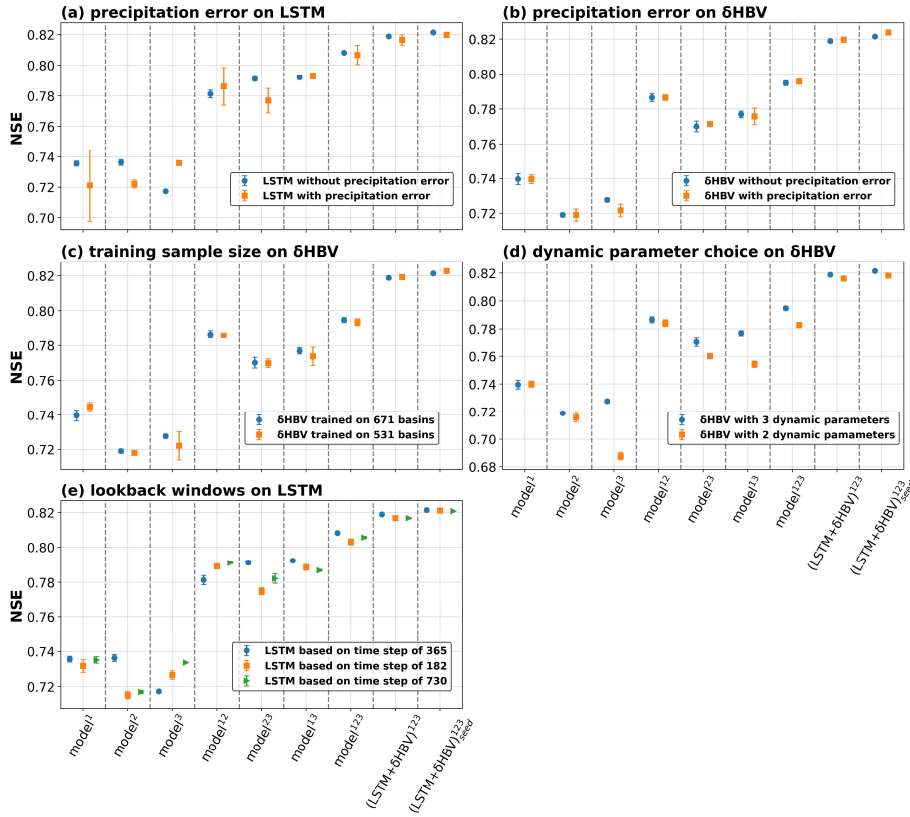


Figure D5. Simulation performance (NSE) under the temporal test: (a) LSTM model with and without a 10% precipitation error ($\text{precipitation} \times 1.1$); (b) δHBV model with and without a 10% precipitation error; (c) δHBV model trained on 671 versus 531 basins; (d) δHBV model with 3 versus 2 dynamic parameters; (e) δHBV model using time steps of 365, 182, and 730 days. Individual and ensemble groups are distinguished along the x-axis. Ensemble benefits are indicated by the gap between columns of the same color within each panel—columns 1–7 correspond to individual LSTM or δHBV groups, and the last two columns correspond to LSTM+ δHBV ensembles.

1043 Table D1. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high
 1044 (highRMSE), and middle (midRMSE) flows based on 531 basins under the temporal test. The
 1045 values are the mean of three simulations run with different random seeds.

Temporal	Number	Daymet	NLDAS	Maurer
LSTM	NSE	0.735639	0.736301	0.717337
	KGE	0.789375	0.782555	0.760575
	RMSE	1.21088	1.19847	1.27723
	PBIAS	4.04818	5.99486	1.58911
	lowRMSE	0.0596913	0.0602381	0.0545577
	highRMSE	2.70508	2.89684	2.97028
	midRMSE	0.196039	0.210022	0.219922
δ HBV	NSE	0.7416417 39688	0.7197767 1903	0.7291427 27669
	KGE	0.7695227 7033	0.7339837 30753	0.7604537 62022
	RMSE	1.1786418 752	1.2686426 239	1.2208923 193
	PBIAS	4.658985 07898	- 0.2289251 4449	3.1474265 263
	lowRMSE	0.0598199 060906	0.0646098 063581	0.0627206 063466
	highRMSE	2.6918684	3.1519513	2.7162968

Formatted: Header

Formatted Table

		<u>79</u>	<u>011</u>	<u>45</u>
	midRMSE	<u>0.2287312</u> <u>26595</u>	<u>0.2450142</u> <u>45242</u>	<u>0.2307252</u> <u>30125</u>
LSTM+δHBV	NSE	<u>0.7878717</u> <u>87545</u>	<u>0.7931687</u> <u>94053</u>	<u>0.7916377</u> <u>90903</u>
	KGE	<u>0.7963227</u> <u>94412</u>	<u>0.7836127</u> <u>8383</u>	<u>0.7847547</u> <u>86854</u>
	RMSE	<u>1.0760407</u> <u>77</u>	<u>1.0674607</u> <u>16</u>	<u>1.0692107</u> <u>141</u>
	PBIAS	<u>4.8257259</u> <u>065</u>	<u>3.0815330</u> <u>53</u>	<u>3.1984145</u> <u>501</u>
	lowRMSE	<u>0.0599687</u> <u>059955</u>	<u>0.0593688</u> <u>059565</u>	<u>0.0541188</u> <u>054838</u>
	highRMSE	<u>2.6966570</u> <u>216</u>	<u>2.8224588</u> <u>511</u>	<u>2.6942569</u> <u>633</u>
	midRMSE	<u>0.2042612</u> <u>0394</u>	<u>0.2184982</u> <u>14726</u>	<u>0.2143252</u> <u>12514</u>

Formatted: Header

1049 Table D1 (continued). Median NSE, KGE, RMSE, PBIAS, and RMSE values under low
 1050 (lowRMSE), high (highRMSE), and middle (midRMSE) flows based on 531 basins under the
 1051 temporal test. The values are the mean of three simulations run with different random seeds.

Temporal	Number	Daymet+NLDAS	Daymet+Maurer	NLDAS+Maurer	A
LSTM	NSE	0.781275	0.791158	0.792144	0.808176
	KGE	0.800955	0.795026	0.794441	0.803476
	RMSE	1.09103	1.06374	1.06701	1.01395
	PBIAS	5.17159	3.34362	4.5305	4.48263
	lowRMSE	0.0636155	0.0582563	0.0566306	0.0613625
	highRMSE	2.70218	2.71366	2.78962	2.67803
	midRMSE	0.194849	0.199809	0.206653	0.197469
δ HBV	NSE	0.786118786562	0.77093977012	0.777651776938	0.794455794796
	KGE	0.772697773732	0.776781778557	0.767756768854	0.77669277834
	RMSE	1.0798408362	1.1267112584	1.1087810875	1.0580806118
	PBIAS	1.8596291507	4.2627828194	1.791342.03584	2.5906371021
	lowRMSE	0.0627661061667	0.0597778060679	0.0623739062765	0.0617863061539
	highRMSE	2.9427493961	2.730547394	2.8758388758	2.8451184994

	midRMSE	0.231984230576	0.219738220743	0.228451230272	0.230136228375
LSTM+δHBV	NSE	0.8108811825	0.810476809964	0.812144811316	0.81866818907
	KGE	0.79586797564	0.796202797635	0.78608878735	0.794257794936
	RMSE	1.016201938	1.0167601755	1.025150314	1.0007700067
	PBIAS	4.1307714594	4.0809623333	3.2645819652	3.897288096
	lowRMSE	0.0599350603	0.0575384058022	0.0558506057882	0.0581869059221
	highRMSE	2.7613375275	2.6864267122	2.7824281393	2.7139270606
	midRMSE	0.208476207637	0.207761205965	0.213433213191	0.208582207905

1053 Table D2. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high
 1054 (highRMSE), and middle (midRMSE) flows based on 531 basins under the PUB test. The values
 1055 are the mean of three simulations run with different random seeds.

PUB	Number	Daymet	NLDAS	Maurer
LSTM	NSE	0.702636	0.695496	0.694156
	KGE	0.693998	0.677438	0.6909
	RMSE	1.31714	1.3394	1.34233
	PBIAS	0.669018	0.283106	0.936582
	lowRMSE	0.087648	0.088393	0.086873
	highRMSE	4.2852	4.49292	4.16042
	midRMSE	0.354458	0.364921	0.368124
δHBV	NSE	0.706809	0.670636	0.682998
	KGE	0.703137	0.66566	0.686912
	RMSE	1.35541	1.41185	1.37942
	PBIAS	1.49234	-2.43395	0.291966
	lowRMSE	0.0798196	0.0808967	0.0846775
	highRMSE	4.21648	4.49582	4.18003
	midRMSE	0.335159	0.351271	0.356903
LSTM+δHBV	NSE	0.74227	0.723778	0.72202

Formatted: Header

Formatted Table

1056
1057
1058

KGE	0.715931	0.690154	0.707292
RMSE	1.24887	1.278	1.26697
PBIAS	1.27863	-0.599778	0.903464
lowRMSE	0.0816748	0.0795686	0.0825691
highRMSE	4.08432	4.23483	3.94929
midRMSE	0.327459	0.33851	0.347169

Formatted: Header

Formatted: Header

1059 Table D2 (continued). Median NSE, KGE, RMSE, PBIAS, and RMSE values under low
 1060 (lowRMSE), high (highRMSE), and middle (midRMSE) flows based on 531 basins under the
 1061 PUB test. The values are the mean of three simulations run with different random seeds.

PUB	Number	Daymet+NLDAS	Daymet+Maurer	NLDAS+Maurer	Formatted Table
LSTM	NSE	0.757853	0.749151	0.753136	0.768181
	KGE	0.713319	0.720099	0.716497	0.727143
	RMSE	1.18251	1.22254	1.19718	1.15026
	PBIAS	0.320396	0.931656	0.766216	0.970047
	lowRMSE	0.0875191	0.0864129	0.0835341	0.0874717
	highRMSE	4.1296	4.06602	4.17217	4.0061
	midRMSE	0.334683	0.349856	0.342819	0.333534
δ HBV	NSE	0.748916	0.734052	0.733955	0.757749
	KGE	0.699768	0.714323	0.69436	0.714048
	RMSE	1.26852	1.27637	1.27244	1.23229
	PBIAS	0.0446112	1.212	-1.04135	0.201809
	lowRMSE	0.0808293	0.0792486	0.0814476	0.0808359
	highRMSE	4.19575	3.97788	4.21623	4.07419
	midRMSE	0.311826	0.33668	0.339257	0.318165
LSTM+ δ HBV	NSE	0.780625	0.764866	0.767761	0.785833

Formatted: Header

KGE	0.719781	0.725373	0.715982	0.723972
RMSE	1.14924	1.17659	1.16881	1.13591
PBIAS	0.186062	0.881644	0.405548	0.565489
lowRMSE	0.0805946	0.0814251	0.0817114	0.0826379
highRMSE	3.97373	3.86834	3.88	3.91692
midRMSE	0.313708	0.324777	0.324089	0.323671

1062

1063

1064 Table D3. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high
 1065 (highRMSE), and middle (midRMSE) flows based on 531 basins under the PUR test. The values
 1066 are the mean of three simulations run with different random seeds.

PUR	Number	Daymet	NLDAS	Maurer
LSTM	NSE	0.578365	0.546217	0.56164
	KGE	0.557788	0.559986	0.567231
	RMSE	1.59111	1.63626	1.5833
	PBIAS	-0.575328	-2.77709	-0.623183
	lowRMSE	0.124837	0.118971	0.118695
	highRMSE	5.42346	5.38886	5.05212
	midRMSE	0.498133	0.498442	0.471744
δHBV	NSE	0.622278	0.592306	0.59161
	KGE	0.638818	0.601338	0.620877
	RMSE	1.57189	1.61191	1.63628
	PBIAS	1.27223	-1.60075	1.62709
	lowRMSE	0.10142	0.102975	0.101075
	highRMSE	5.07706	5.16093	4.99602
	midRMSE	0.447879	0.474516	0.439697
LSTM+δHBV	NSE	0.644398	0.618255	0.635444

Formatted: Header

Formatted Table

1067

1068

KGE	0.627481	0.605237	0.615883
RMSE	1.46185	1.5153	1.48393
PBIAS	-0.269697	-0.719505	0.197859
lowRMSE	0.105146	0.100944	0.106272
highRMSE	4.95749	4.99478	4.78638
midRMSE	0.431456	0.4575	0.426126

Formatted: Header

Formatted: Header

1069 Table D3 (continued). Median NSE, KGE, RMSE, PBIAS, and RMSE values under low
 1070 (lowRMSE), high (highRMSE), and middle (midRMSE) flows based on 531 basins under the
 1071 PUR test. The values are the mean of three simulations run with different random seeds.

PUR	Number	Daymet+NLDAS	Daymet+Maurer	NLDAS+Maurer	Formatted Table
LSTM	NSE	0.634398	0.636369	0.626939	0.656228
	KGE	0.59844	0.600371	0.605007	0.612858
	RMSE	1.4434	1.43416	1.43009	1.38042
	PBIAS	-0.547128	-0.687947	-0.865748	-0.543918
	lowRMSE	0.118989	0.120228	0.115004	0.117728
	highRMSE	5.03277	5.02434	4.84415	4.74281
	midRMSE	0.462923	0.455257	0.453912	0.449598
δ HBV	NSE	0.672839	0.644732	0.661231	0.684685
	KGE	0.653841	0.65646	0.6515	0.66205
	RMSE	1.43224	1.50803	1.48604	1.43376
	PBIAS	0.564363	1.55134	-0.156553	0.956961
	lowRMSE	0.0975783	0.0984076	0.100773	0.100807
	highRMSE	4.83843	4.81176	4.72529	4.71255
	midRMSE	0.447828	0.431252	0.433688	0.432018
LSTM+ δ HBV	NSE	0.685032	0.680872	0.679321	0.700814

Formatted: Header

KGE	0.638788	0.647826	0.646782	0.649999
RMSE	1.35303	1.3873	1.36795	1.3185
PBIAS	-0.0150729	0.406127	-0.135091	-0.0232668
lowRMSE	0.103284	0.101814	0.104528	0.102916
highRMSE	4.80178	4.72583	4.70024	4.70713
midRMSE	0.426819	0.411727	0.41573	0.41081

1072

1073

1074 Table D4. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high
1075 (highRMSE), and middle (midRMSE) flows based on 531 basins under the temporal, PUB, and
1076 PUR tests of $LSTM^{multi}$, $(LSTM + \delta HBV)^{123} + LSTM^{multi}$, their “seed” version, and
1077 $(LSTM + \delta HBV)_{seed}^{123}$.

Test	Metric	$LSTM^{multi}$	$(LSTM + \delta HBV)^{123}$ + $LSTM^{multi}$
Temporal	NSE	0.797448	0.82321
	KGE	0.811064	0.810248
	RMSE	1.05987	0.983168
	PBIAS	3.95241	4.08594
	lowRMSE	0.056221	0.05702
	highRMSE	2.7089	2.58881
	midRMSE	0.183526	0.192442
PUB	NSE	0.750605	0.782727
	KGE	0.71469	0.734731
	RMSE	1.20586	1.11509
	PBIAS	0.475674	0.706777
	lowRMSE	0.0861127	0.0836
	highRMSE	4.13615	3.83009
	midRMSE	0.347562	0.326814

Formatted: Header

Formatted Table

|

|

|

|

|

|

|

|

1078

PUR	NSE	0.623755	0.68923
	KGE	0.593757	0.633971
	RMSE	1.47379	1.31221
	PBIAS	-2.6737	-1.38119
	lowRMSE	0.112434	0.107646
	highRMSE	4.98202	4.59232
	midRMSE	0.501807	0.436811

Formatted: Header

Formatted: Header

1079 Table D4 (continued). Median NSE, KGE, RMSE, PBIAS, and RMSE values under low
 1080 (lowRMSE), high (highRMSE), and middle (midRMSE) flows based on 531 basins under the
 1081 temporal, PUB, and PUR tests of $LSTM^{multi}$, $(LSTM + \delta HBV)^{123}$ + $LSTM^{multi}$, their “seed”
 1082 version, and $(LSTM + \delta HBV)^{123}_{seed}$.

Formatted Table

Test	Metric	$(LSTM + \delta HBV)^{123}_{seed}$	$LSTM^{multi}_{seed}$	$(LSTM + \delta HBV)^{123}_{seed}$ + $LSTM^{multi}_{seed}$
Temporal	NSE	0.821444	0.81992	0.829385
	KGE	0.795317	0.82078	0.812581
	RMSE	0.99455	1.00908	0.967779
	PBIAS	3.99009	4.09469	4.08882
	lowRMSE	0.059782	0.057346	0.057015
	highRMSE	2.7279	2.62815	2.58384
	midRMSE	0.209943	0.183656	0.195557
PUB	NSE	0.793673	0.781175	0.790921
	KGE	0.726188	0.736191	0.739284
	RMSE	1.12957	1.13079	1.09176
	PBIAS	0.370674	1.13671	0.869057
	lowRMSE	0.083423	0.084038	0.085728
	highRMSE	3.89363	3.93473	3.79505
	midRMSE	0.323045	0.329772	0.325627

PUR	NSE	0.705154	0.665723	0.701504
	KGE	0.651538	0.614649	0.64373
	RMSE	1.30377	1.3727	1.2851
	PBIAS	-0.283645	-2.74069	-1.39149
	lowRMSE	0.100525	0.111229	0.108121
	highRMSE	4.74889	4.88127	4.58344
	midRMSE	0.406797	0.473783	0.432447

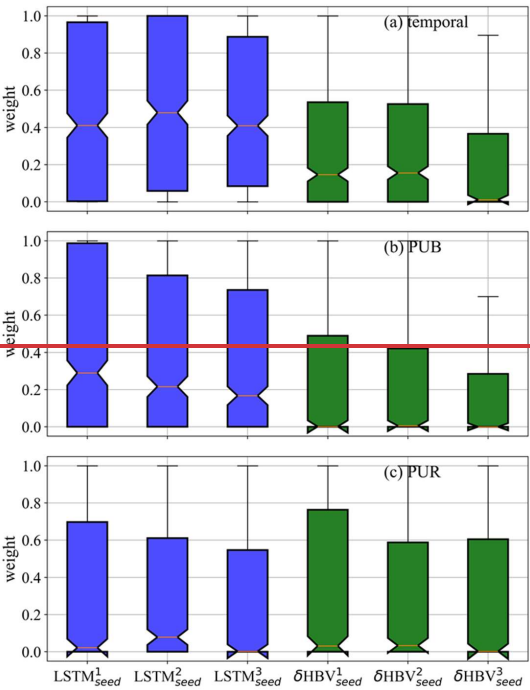
Formatted: Header

Formatted: Indent: First line: 0.33"

1084 Table D5. Median NSE values based on ten different random seeds during the temporal test.
1085 Each number (1 through 10) represents metric values calculated for an individual simulation
1086 based on only one random seed. “Seed” indicates metric values calculated by averages of these
1087 ten simulations based on different random seeds, while “mean” denotes the average of metrics
1088 from 1-10 individual simulations (visualized in Figure C1).

Number	$LSTM^{multi}$	$(LSTM + \delta HBV)^{123}$	$(LSTM + \delta HBV)^{123} + LSTM^{multi}$
1	0.797742	0.818436	0.82315
2	0.795312	0.820188	0.823559
3	0.799291	0.818097	0.822922
4	0.796388	0.818251	0.821791
5	0.791192	0.818285	0.820132
6	0.795691	0.81966	0.823268
7	0.795912	0.821511	0.82352
8	0.796625	0.81831	0.825204
9	0.794062	0.804959	0.816497
10	0.796066	0.817122	0.82169
Seed	0.82425	0.822528	0.832197
Mean	0.795828	0.817482	0.822173

1090 **Appendix E: Intuitive visualization of the relative contributions of ensemble members**
1091 **based on optimized weights**



Formatted: Header

Formatted: Font color: Auto

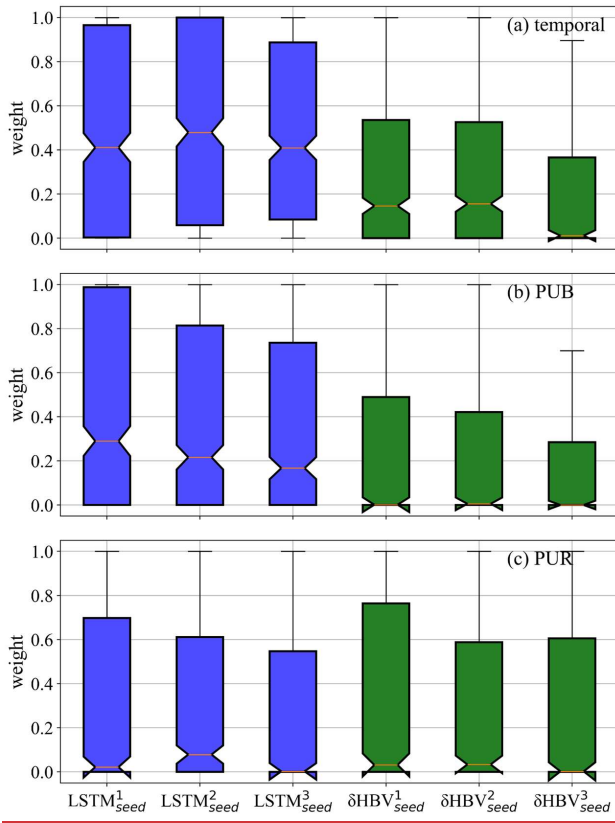


Figure E1. Weights of six components across 531 basins, estimated basin-by-basin using a genetic algorithm based on streamflow observations during the test periods. The weights are normalized by the maximum weight within each ensemble group. These weights are used exclusively for qualitatively analyzing the relative contributions of different ensemble members, with higher values indicating larger relative contributions.

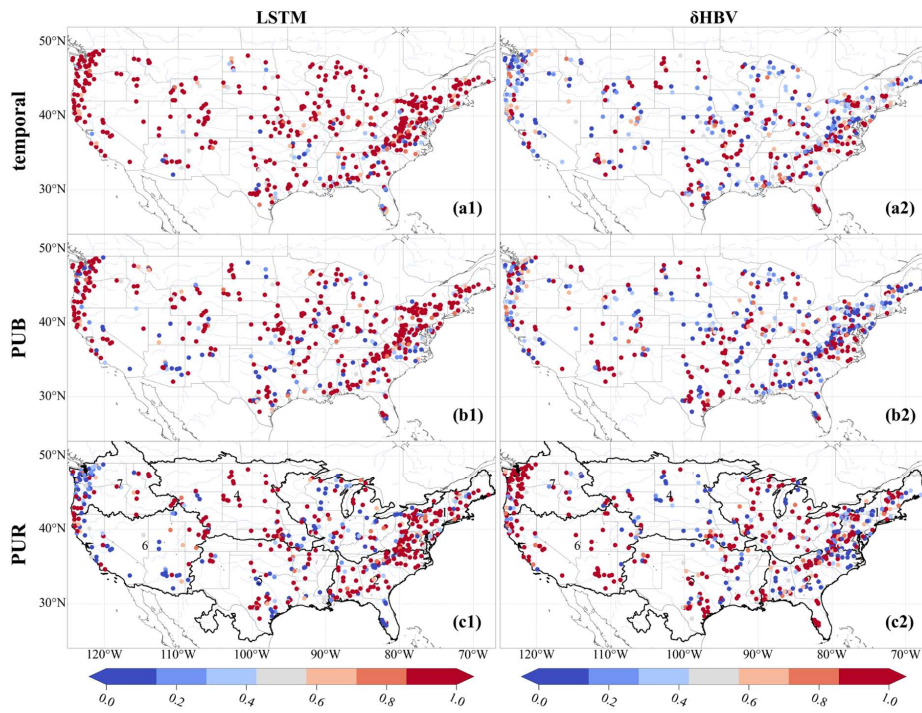
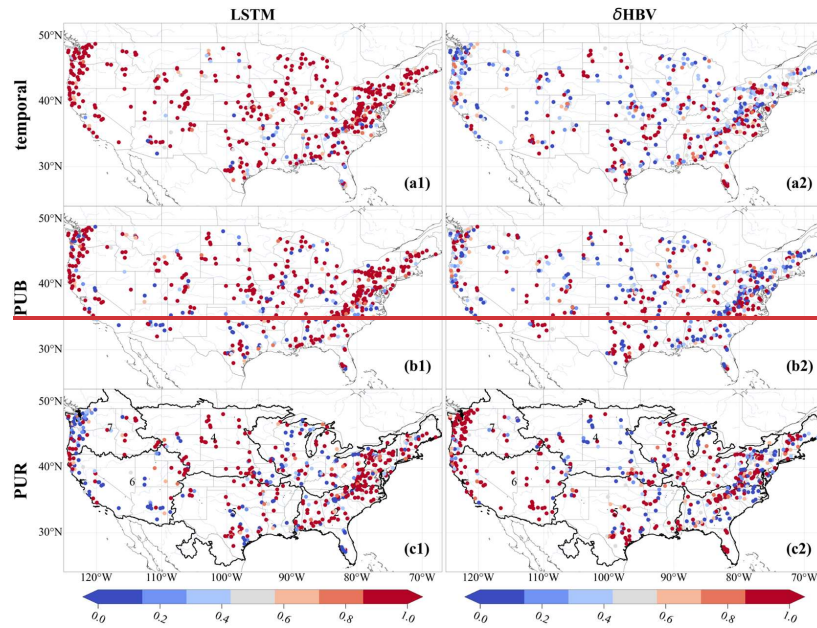


Figure E2. Spatial distributions of weights of the LSTM and δ HBV models, estimated by a

genetic algorithm based on streamflow observations during the test periods. The weights are normalized by the maximum weight within each ensemble group. These weights are used exclusively for qualitatively analyzing the relative contributions of different ensemble members, with higher values indicating larger relative contributions.

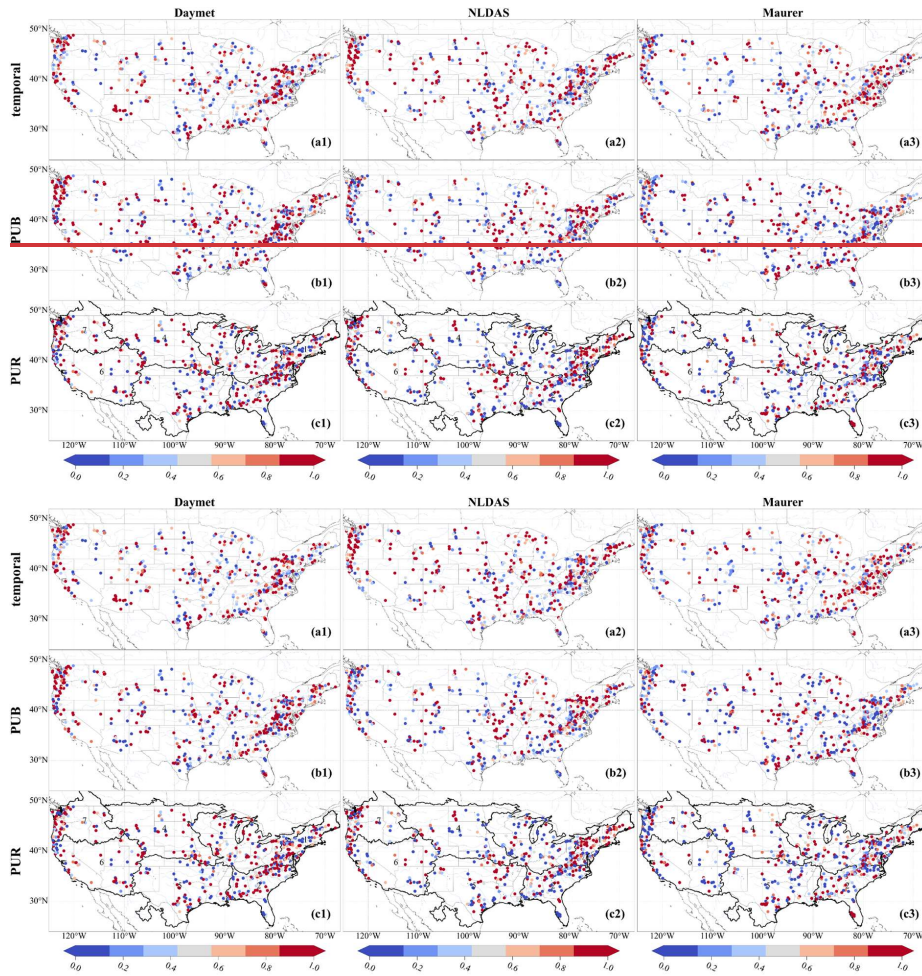


Figure E3. Spatial distributions of weights of the Daymet, NLDAS, and Maurer meteorological forcing datasets, estimated by a genetic algorithm based on streamflow observations during the test periods. The weights are normalized by the maximum weight within each ensemble group. These weights are used exclusively for qualitatively analyzing the relative contributions

|
1115 *of different ensemble members, with higher values indicating larger relative contributions.*
1116
1117
1118

Table E1. Comparisons of metric values between averaged ensemble simulations and optimized weighted simulations, estimated using a genetic algorithm based on streamflow observations during the test periods. The results highlight the potential for further improvements in ensemble simulations.

	Temporal	Averaged	Optimized weighted
Temporal	NSE	0.82144391282144	0.844303212
	KGE	0.795317495795317	0.829996445
	RMSE	0.99455008299455	0.920954559
	PBIAS	3.99009459199009	3.252278013
	lowRMSE	0.059781616059782	0.057137161
	highRMSE	2.727901337279	2.451194907
	midRMSE	0.20994263209943	0.183127162
PUB	NSE	0.793673	0.842396015
	KGE	0.726188	0.79571295
	RMSE	1.12957	0.987170488
	PBIAS	0.370674	1.023040859
	lowRMSE	0.0834234	0.079807878
	highRMSE	3.89363	3.030715903
	midRMSE	0.323045	0.285110115
PUR	NSE	0.705154	0.790796063
	KGE	0.651538	0.746396324

Formatted: Header

Formatted: Justified, Line spacing: single

Formatted Table

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

	RMSE	1.30377	1.13058149
	PBIAS	-0.283645	0.273698787
	lowRMSE	0.100525	0.093595304
	highRMSE	4.74889	3.665495069
	midRMSE	0.406797	0.351694421

Formatted: Header

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

Formatted: Justified, Line spacing: single

|

1|25

1|26

Formatted: Header

1127 **Code and data availability**

1128 The source codes and datasets utilized in this study are publicly accessible through the

1129 following repositories: The δ HBV modeling framework, including all computational scripts

1130 and documentation, is hosted on Zenodo (<https://doi.org/10.5281/zenodo.7091334>) (Feng et al.,

1131 2023a), with an updated version and comprehensive software release scheduled upon

1132 manuscript acceptance. The implementation of the LSTM architecture is accessible through

1133 Zenodo (<https://doi.org/10.5281/zenodo.6326394>) (Kratzert et al., 2022). The CAMELS

1134 hydrometeorological dataset, which provides the foundational basin characteristics and time

1135 series data used in our analysis, can be obtained via <https://dx.doi.org/10.5065/D6MW2F4D>

1136 (Addor et al., 2017; Newman and Clark, 2014). The streamflow simulations produced in this

1137 study will be made available on Zenodo upon acceptance of the manuscript.

Formatted: Header

Formatted: Font color: Auto

1139 **Author contributions**

1140 PL and CS designed the experiments and PL carried them out. YS developed the modified

1141 δ HBV code. PL prepared the manuscript with contributions from all co-authors.

Formatted: Font color: Auto

1143 **Competing interests**

1144 Chaopeng Shen and Kathryn Lawson have financial interests in HydroSapient, Inc., a

1145 company that could potentially benefit from the results of this research. This interest has been

1146 reviewed by the Pennsylvania State University in accordance with its individual conflict of

1147 interest policy for the purpose of maintaining the objectivity and the integrity of research. The

1148 other authors have no competing interests to declare.

Formatted: Font color: Auto

1150 **Acknowledgments**

1151 PL, CS, and KL were supported by the Office of Biological and Environmental Research

Formatted: Font color: Auto

of the U.S. Department of Energy (contract no. DESC0016605). PJ and MP were also partially supported by California Department of Water Resources Atmospheric River Program Phase III (Grant 4600014294). YS and CS were partially supported by subaward A23-0252-S002 from the Cooperative Institute for Research to Operations in Hydrology (CIROH) through the National Oceanic and Atmospheric Administration (NOAA) Cooperative Agreement (Grant no. NA22NWS4320003).

References

Formatted: Font color: Auto

- Aboelyazeed, D., Xu, C., Hoffman, F. M., Liu, J., Jones, A. W., Rackauckas, C., Lawson, K., and Shen, C.: A differentiable, physics-informed ecosystem modeling and learning framework for large-scale inverse problems: demonstration with photosynthesis simulations, *Biogeosciences*, 20, 2671–2692, <https://doi.org/10.5194/bg-20-2671-2023>, 2023.
- Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-sample studies, *Hydrol. Earth Syst. Sci.*, 21, 5293–5313, <https://doi.org/10.5194/hess-21-5293-2017>, 2017.
- Aghakouchak, A. and Habib, E.: Application of a Conceptual Hydrologic Model in Teaching Hydrologic Processes, *International Journal of Engineering Education*, 26, 2010.
- Bandai, T. and Ghezzehei, T. A.: Physics-informed neural networks with monotonicity constraints for Richardson-Richards equation: Estimation of constitutive relationships and soil water flux density from volumetric water content measurements, *Water Resources Research*, 57, e2020WR027642, <https://doi.org/10.1029/2020wr027642>, 2021.
- Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Dutra, E., Fink, G., Orth, R., and Schellekens, J.: Global evaluation of runoff from 10 state-of-the-art hydrological models, *Hydrology and Earth System Sciences*, 21, 2881–2903, <https://doi.org/10.5194/hess-21-2881-2017>, 2017.
- Beck, H. E., Pan, M., Lin, P., Seibert, J., Dijk, A. I. J. M. van, and Wood, E. F.: Global fully distributed parameter regionalization based on observed streamflow from 4,229 headwater catchments, *Journal of Geophysical Research: Atmospheres*, 125, e2019JD031485, <https://doi.org/10.1029/2019JD031485>, 2020.
- Behnke, R., Vavrus, S., Allstadt, A., Albright, T., Thogmartin, W. E., and Radeloff, V. C.: Evaluation of downscaled, gridded climate data for the conterminous United States, *Ecological Applications*, 26, 1338–1351, <https://doi.org/10.1002/15-1061>, 2016.
- Bell, V. A. and Moore, R. J.: The sensitivity of catchment runoff models to rainfall data at different spatial scales, *Hydrology and Earth System Sciences*, 4, 653–667, <https://doi.org/10.5194/hess-4-653-2000>, 2000.
- Bergström, S.: Development and application of a conceptual runoff model for Scandinavian

- catchments, PhD Thesis, Swedish Meteorological and Hydrological Institute (SMHI), Norköping, Sweden, 1976.
- Bergström, S.: The HBV model—its structure and applications, SMHI, 1992.
- Bindas, T., Tsai, W.-P., Liu, J., Rahmani, F., Feng, D., Bian, Y., Lawson, K., and Shen, C.: Improving river routing using a differentiable Muskingum-Cunge model and physics-informed machine learning, *Water Resources Research*, 60, e2023WR035337, <https://doi.org/10.1029/2023WR035337>, 2024.
- Bodnar, C., Bruinsma, W. P., Lucic, A., Stanley, M., Allen, A., Brandstetter, J., Garvan, P., Riechert, M., Weyn, J. A., Dong, H., Gupta, J. K., Thambiratnam, K., Archibald, A. T., Wu, C.-C., Heider, E., Welling, M., Turner, R. E., and Perdikaris, P.: A foundation model for the Earth system, *Nature*, 641, 1180–1187, <https://doi.org/10.1038/s41586-025-09005-y>, 2025.
- Brunner, M. I., Slater, L., Tallaksen, L. M., and Clark, M.: Challenges in modeling and predicting floods and droughts: A review, *WIREs Water*, 8, e1520, <https://doi.org/10.1002/wat2.1520>, 2021.
- Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T., and Hay, L. E.: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between hydrological models, *Water Resources Research*, 44, <https://doi.org/10/chvc6k>, 2008.
- Clark, M. P., Nijssen, B., Lundquist, J. D., Kavetski, D., Rupp, D. E., Woods, R. A., Freer, J. E., Gutmann, E. D., Wood, A. W., Brekke, L. D., Arnold, J. R., Gochis, D. J., and Rasmussen, R. M.: A unified approach for process-based hydrologic modeling: 1. Modeling concept, *Water Resources Research*, 51, 2498–2514, <https://doi.org/10/f7db99>, 2015.
- Clark, M. P., Wilby, R. L., Gutmann, E. D., Vano, J. A., Gangopadhyay, S., Wood, A. W., Fowler, H. J., Prudhomme, C., Arnold, J. R., and Brekke, L. D.: Characterizing uncertainty of the hydrologic impacts of climate change, *Curr Clim Change Rep*, 2, 55–64, <https://doi.org/10.1007/s40641-016-0034-x>, 2016.
- Dion, P., Martel, J.-L., and Arsenault, R.: Hydrological ensemble forecasting using a multi-model framework, *Journal of Hydrology*, 600, 126537, <https://doi.org/10.1016/j.jhydrol.2021.126537>, 2021.
- Feng, D., Fang, K., and Shen, C.: Enhancing streamflow forecast and extracting insights using long-short term memory networks with data integration at continental scales, *Water Resources Research*, 56, e2019WR026793, <https://doi.org/10.1029/2019WR026793>, 2020.
- Feng, D., Lawson, K., and Shen, C.: Mitigating prediction error of deep learning streamflow models in large data-sparse regions with ensemble modeling and soft data, *Geophysical Research Letters*, 48, e2021GL092999, <https://doi.org/10.1029/2021GL092999>, 2021.
- Feng, D., Liu, J., Lawson, K., and Shen, C.: Differentiable, learnable, regionalized process-based models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy, *Water Resources Research*, 58, e2022WR032404, <https://doi.org/10.1029/2022WR032404>, 2022.
- Feng, D., Shen, C., Liu, J., Lawson, K., and Beck, H.: differentiable parameter learning (dPL)

- + HBV hydrologic model, <https://doi.org/10.5281/zenodo.7943626>, 2023a.
- Feng, D., Beck, H., Lawson, K., and Shen, C.: The suitability of differentiable, physics-informed machine learning hydrologic models for ungauged regions and climate change impact assessment, *Hydrology and Earth System Sciences*, 27, 2357–2373, <https://doi.org/10.5194/hess-27-2357-2023>, 2023b.
- Frame, J. M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., Qualls, L. M., Gupta, H. V., and Nearing, G. S.: Deep learning rainfall–runoff predictions of extreme events, *Hydrology and Earth System Sciences*, 26, 3377–3392, <https://doi.org/10.5194/hess-26-3377-2022>, 2022.
- Hanazaki, R., Yamazaki, D., and Yoshimura, K.: Development of a reservoir flood control scheme for global flood models, *JAMES*, 14, e2021MS002944, <https://doi.org/10.1029/2021MS002944>, 2022.
- Hargreaves, G. H.: Defining and using reference evapotranspiration, *Journal of Irrigation and Drainage Engineering*, 120, 1132–1139, [https://doi.org/10.1061/\(ASCE\)0733-9437\(1994\)120:6\(1132\)](https://doi.org/10.1061/(ASCE)0733-9437(1994)120:6(1132)), 1994.
- He, Y., Chen, M., Wen, Y., Duan, Q., Yue, S., Zhang, J., Li, W., Sun, R., Zhang, Z., Tao, R., Tang, W., and Lü, G.: An open online simulation strategy for hydrological ensemble forecasting, *Environmental Modelling & Software*, 174, 105975, <https://doi.org/10.1016/j.envsoft.2024.105975>, 2024.
- Heidari, H., Arabi, M., Warziniack, T., and Kao, S.-C.: Assessing shifts in regional hydroclimatic conditions of U.S. river basins in response to climate change over the 21st century, *Earth's Future*, 8, e2020EF001657, <https://doi.org/10.1029/2020EF001657>, 2020.
- Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, *Neural Computation*, 9, 1735–1780, <https://doi.org/10.1162/neco.1997.9.8.1735>, 1997.
- Jiang, S., Zheng, Y., and Solomatine, D.: Improving AI system awareness of geoscience knowledge: Symbiotic integration of physical approaches and deep learning, *Geophys. Res. Lett.*, 47, e2020GL088229, <https://doi.org/10.1029/2020GL088229>, 2020.
- Kling, H., Fuchs, M., and Paulin, M.: Runoff conditions in the upper Danube basin under an ensemble of climate change scenarios, *Journal of Hydrology*, 424–425, 264–277, <https://doi.org/10.1016/j.jhydrol.2012.01.011>, 2012.
- Kraft, B., Jung, M., Körner, M., Koirala, S., and Reichstein, M.: Towards hybrid modeling of the global hydrological cycle, *Hydrology and Earth System Sciences*, 26, 1579–1614, <https://doi.org/10.5194/hess-26-1579-2022>, 2022.
- Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall-Runoff modelling using Long-Short-Term-Memory (LSTM) networks, *Hydrology and Earth System Sciences*, 22, 6005–6022, <https://doi.org/10.17605/OSF.IO/QV5JZ>, 2018.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward improved predictions in ungauged basins: Exploiting the power of machine learning, *Water Resources Research*, 55, 11344–11354, <https://doi.org/10.1029/2019WR025468>, 2019.

- Kratzert, F., Klotz, D., Hochreiter, S., and Nearing, G. S.: A note on leveraging synergy in multiple meteorological data sets with deep learning for rainfall–runoff modeling, *Hydrology and Earth System Sciences*, 25, 2685–2703, <https://doi.org/10.5194/hess-25-2685-2021>, 2021.
- Kratzert, F., Gauch, M., Nearing, G., and Klotz, D.: NeuralHydrology — A Python library for Deep Learning research in hydrology, <https://doi.org/10.5281/zenodo.6326394>, 2022.
- Leube, P. C., de Barros, F. P. J., Nowak, W., and Rajagopal, R.: Towards optimal allocation of computer resources: Trade-offs between uncertainty quantification, discretization and model reduction, *Environmental Modelling & Software*, 50, 97–107, <https://doi.org/10.1016/j.envsoft.2013.08.008>, 2013.
- Li, P., Zha, Y., Shi, L., Tso, C. H. M., Zhang, Y., and Zeng, W.: Comparison of the use of a physical-based model with data assimilation and machine learning methods for simulating soil water dynamics, *Journal of Hydrology*, 584, 124692, <https://doi.org/10.1016/j.jhydrol.2020.124692>, 2020a.
- Li, P., Zha, Y., Tso, C. H. M., Shi, L., Yu, D., Zhang, Y., and Zeng, W.: Data assimilation of uncalibrated soil moisture measurements from frequency-domain reflectometry, *Geoderma*, 374, 114432, <https://doi.org/10.1016/j.geoderma.2020.114432>, 2020b.
- Li, P., Zha, Y., Shi, L., and Zhong, H.: Identification of the terrestrial water storage change features in the North China Plain via independent component analysis, *Journal of Hydrology: Regional Studies*, 38, 100955, <https://doi.org/10.1016/j.ejrh.2021.100955>, 2021.
- Li, P., Zha, Y., Zuo, B., and Zhang, Y.: A family of soil water retention models based on sigmoid functions, *Water Resources Research*, 59, e2022WR033160, <https://doi.org/10.1029/2022WR033160>, 2023a.
- Li, P., Zha, Y., and Tso, C.-H. M.: Reconstructing GRACE-derived terrestrial water storage anomalies with in-situ groundwater level measurements and meteorological forcing data, *Journal of Hydrology: Regional Studies*, 50, 101528, <https://doi.org/10.1016/j.ejrh.2023.101528>, 2023b.
- Li, P., Zha, Y., Zhang, Y., Michael Tso, C.-H., Attinger, S., Samaniego, L., and Peng, J.: Deep learning integrating scale conversion and pedo-transfer function to avoid potential errors in cross-scale transfer, *Water Resources Research*, 60, e2023WR035543, <https://doi.org/10.1029/2023WR035543>, 2024.
- Lin, Y., Wang, D., Zhu, J., Sun, W., Shen, C., and Shangguan, W.: Development of objective function-based ensemble model for streamflow forecasts, *Journal of Hydrology*, 632, 130861, <https://doi.org/10.1016/j.jhydrol.2024.130861>, 2024.
- Lins, H. F. and Slack, J. R.: Streamflow trends in the United States, *Geophysical Research Letters*, 26, 227–230, <https://doi.org/10/d5zbbd>, 1999.
- Liu, J., Rahmani, F., Lawson, K., and Shen, C.: A multiscale deep learning model for soil moisture integrating satellite and in situ data, *Geophysical Research Letters*, 49, e2021GL096847, <https://doi.org/10.1029/2021GL096847>, 2022.
- Liu, J., Bian, Y., Lawson, K., and Shen, C.: Probing the limit of hydrologic predictability with the Transformer network, *Journal of Hydrology*, 637, 131389, <https://doi.org/10.1016/j.jhydrol.2024.131389>, 2024.

- <https://doi.org/10.1016/j.jhydrol.2024.131389>, 2024.
- Mai, J., Craig, J. R., Tolson, B. A., and Arsenaault, R.: The sensitivity of simulated streamflow to individual hydrologic processes across North America, *Nat Commun*, 13, 455, <https://doi.org/10.1038/s41467-022-28010-7>, 2022.
- Maurer, E. P., Wood, A. W., Adam, J. C., Lettenmaier, D. P., and Nijssen, B.: A long-term hydrologically based dataset of land surface fluxes and states for the conterminous United States, *Journal of Climate*, 15, 3237–3251, [https://doi.org/10.1175/1520-0442\(2002\)015<3237:ALTHBD>2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015<3237:ALTHBD>2.0.CO;2), 2002.
- Moges, E., Demissie, Y., and Li, H.-Y.: Hierarchical mixture of experts and diagnostic modeling approach to reduce hydrologic model structural uncertainty, *Water Resources Research*, 52, 2551–2570, <https://doi.org/10.1002/2015WR018266>, 2016.
- Nai, C., Liu, X., Tang, Q., Liu, L., Sun, S., and Gaffney, P. P. J.: A novel strategy for automatic selection of cross-basin data to improve local machine learning-based runoff models, *Water Resources Research*, 60, e2023WR035051, <https://doi.org/10.1029/2023WR035051>, 2024.
- Narkhede, M. V., Bartakke, P. P., and Sutaone, M. S.: A review on weight initialization strategies for neural networks, *Artificial Intelligence Review*, 55, 291–322, <https://doi.org/10.1007/s10462-021-10033-z>, 2022.
- Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I — A discussion of principles, *Journal of Hydrology*, 10, 282–290, [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6), 1970.
- Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz, D., Kratzert, F., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev, G., Shenzis, S., Tekalign, T. Y., Weitzner, D., and Matias, Y.: Global prediction of extreme floods in ungauged watersheds, *Nature*, 627, 559–563, <https://doi.org/10.1038/s41586-024-07145-1>, 2024.
- Newman, A. J. and Clark, M.: A large-sample watershed-scale hydrometeorological dataset for the contiguous USA, <https://doi.org/10.5065/D6MW2F4D>, 2014.
- Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., Nearing, G., Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., and Nearing, G.: Benchmarking of a Physically Based Hydrologic Model, *Journal of Hydrometeorology*, 18, 2215–2225, <https://doi.org/10/gbwr9s>, 2017.
- Newman, A. J., Clark, M. P., Longman, R. J., and Giambelluca, T. W.: Methodological intercomparisons of station-based gridded meteorological products: Utility, limitations, and paths forward, <https://doi.org/10.1175/JHM-D-18-0114.1>, 2019.
- Ouyang, W., Lawson, K., Feng, D., Ye, L., Zhang, C., and Shen, C.: Continental-scale streamflow modeling of basins with reservoirs: Towards a coherent deep-learning-based strategy, *Journal of Hydrology*, 599, 126455, <https://doi.org/10.1016/j.jhydrol.2021.126455>, 2021.
- Paul, P. K., Zhang, Y., Ma, N., Mishra, A., Panigrahy, N., and Singh, R.: Selecting hydrological models for developing countries: Perspective of global, continental, and country scale models over catchment scale models, *Journal of Hydrology*, 600, 126561, <https://doi.org/10.1016/j.jhydrol.2021.126561>, 2021.

- <https://doi.org/10.1016/j.jhydrol.2021.126561>, 2021.
- Rahmani, F., Appling, A., Feng, D., Lawson, K., and Shen, C.: Identifying structural priors in a hybrid differentiable model for stream water temperature modeling, *Water Resources Research*, 59, e2023WR034420, <https://doi.org/10.1029/2023WR034420>, 2023.
- Reichle, R. H. and Koster, R. D.: Assessing the impact of horizontal error correlations in background fields on soil moisture estimation, *Journal of Hydrometeorology*, 4, 1229–1242, [https://doi.org/10.1175/1525-7541\(2003\)004<1229:ATIOHE>2.0.CO;2](https://doi.org/10.1175/1525-7541(2003)004<1229:ATIOHE>2.0.CO;2), 2003.
- Ryan Bellmore, J., Duda, J. J., Craig, L. S., Greene, S. L., Torgersen, C. E., Collins, M. J., and Vittum, K.: Status and trends of dam removal research in the United States, *WIREs Water*, 4, e1164, <https://doi.org/10.1002/wat2.1164>, 2017.
- Sawadekar, K., Song, Y., Pan, M., Beck, H., McCrary, R., Ullrich, P., Lawson, K., and Shen, C.: Improving differentiable hydrologic modeling with interpretable forcing fusion, *J. Hydrol.*, 659, 133320, <https://doi.org/10.1016/j.jhydrol.2025.133320>, 2025.
- Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., Baity-Jesi, M., Fenicia, F., Kifer, D., Li, L., Liu, X., Ren, W., Zheng, Y., Harman, C. J., Clark, M., Farthing, M., Feng, D., Kumar, P., Aboelyazeed, D., Rahmani, F., Song, Y., Beck, H. E., Bindas, T., Dwivedi, D., Fang, K., Höge, M., Rackauckas, C., Mohanty, B., Roy, T., Xu, C., and Lawson, K.: Differentiable modelling to unify machine learning and physical models for geosciences, *Nat Rev Earth Environ*, 4, 552–567, <https://doi.org/10.1038/s43017-023-00450-9>, 2023.
- Solanki, H., Vegad, U., Kushwaha, A., and Mishra, V.: Improving streamflow prediction using multiple hydrological models and machine learning methods, *Water Resources Research*, 61, e2024WR038192, <https://doi.org/10.1029/2024WR038192>, 2025.
- Song, Y., Bindas, T., Shen, C., Ji, H., Knoben, W. J. M., Lonzarich, L., Clark, M. P., Liu, J., van Werkhoven, K., Lemont, S., Denno, M., Pan, M., Yang, Y., Rapp, J., Kumar, M., Rahmani, F., Thébault, C., Sawadekar, K., and Lawson, K.: High-resolution national-scale water modeling is enhanced by multiscale differentiable physics-informed machine learning, <https://doi.org/10.22541/essoar.172736277.74497104/v1>, 26 September 2024a.
- Song, Y., Knoben, W. J. M., Clark, M. P., Feng, D., Lawson, K., Sawadekar, K., and Shen, C.: When ancient numerical demons meet physics-informed machine learning: adjoint-based gradients for implicit differentiable modeling, *Hydrology and Earth System Sciences*, 28, 3051–3077, <https://doi.org/10.5194/hess-28-3051-2024>, 2024b.
- Song, Y., Bindas, T., Shen, C., Ji, H., Knoben, W. J. M., Lonzarich, L., Clark, M. P., Liu, J., van Werkhoven, K., Lamont, S., Denno, M., Pan, M., Yang, Y., Rapp, J., Kumar, M., Rahmani, F., Thébault, C., Adkins, R., Halgren, J., Patel, T., Patel, A., Sawadekar, K. A., and Lawson, K.: High-resolution national-scale water modeling is enhanced by multiscale differentiable physics-informed machine learning, *Water Resour. Res.*, 61, e2024WR038928, <https://doi.org/10.1029/2024WR038928>, 2025a.
- Song, Y., Sawadekar, K., Frame, J. M., Pan, M., Clark, M., Knoben, W. J. M., Wood, A. W., Lawson, K. E., Patel, T., and Shen, C.: Physics-informed, differentiable hydrologic models for capturing unseen extreme events, <https://doi.org/10.22541/essoar.172304428.82707157/v2>, 2025b.

- Thornton, P. E., Running, S. W., and White, M. A.: Generating surfaces of daily meteorological variables over large regions of complex terrain, *Journal of Hydrology*, 190, 214–251, [https://doi.org/10.1016/S0022-1694\(96\)03128-9](https://doi.org/10.1016/S0022-1694(96)03128-9), 1997.
- Tsai, W.-P., Feng, D., Pan, M., Beck, H., Lawson, K., Yang, Y., Liu, J., and Shen, C.: From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling, *Nat Commun*, 12, 5988, <https://doi.org/10.1038/s41467-021-26107-z>, 2021.
- Wada, Y., de Graaf, I. E. M., and van Beek, L. P. H.: High-resolution modeling of human and climate impacts on global water resources, *Journal of Advances in Modeling Earth Systems*, 8, 735–763, <https://doi.org/10/f8wgpv>, 2016.
- Wang, N., Zhang, D., Chang, H., and Li, H.: Deep learning of subsurface flow via theory-guided neural network, *Journal of Hydrology*, 584, 124700, <https://doi.org/10.1016/j.jhydrol.2020.124700>, 2020.
- West, B. D., Maxwell, R. M., and Condon, L. E.: A scalable and modular reservoir implementation for large-scale integrated hydrologic simulations, *Hydrology and Earth System Sciences*, 29, 245–259, <https://doi.org/10.5194/hess-29-245-2025>, 2025.
- Wilbrand, K., Taormina, R., ten Veldhuis, M.-C., Visser, M., Hrachowitz, M., Nuttall, J., and Dahm, R.: Predicting streamflow with LSTM networks using global datasets, *Front. Water*, 5, <https://doi.org/10.3389/frwa.2023.1166124>, 2023.
- Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H., Meng, J., Livneh, B., Lettenmaier, D., Koren, V., Duan, Q., Mo, K., Fan, Y., and Mocko, D.: Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of model products, *Journal of Geophysical Research: Atmospheres*, 117, <https://doi.org/10.1029/2011JD016048>, 2012.
- Xie, K., Liu, P., Zhang, J., Han, D., Wang, G., and Shen, C.: Physics-guided deep learning for rainfall-runoff modeling by considering extreme events and monotonic relationships, *Journal of Hydrology*, 603, 127043, <https://doi.org/10.1016/j.jhydrol.2021.127043>, 2021.
- Yao, L., Libera, D. A., Kheimi, M., Sankarasubramanian, A., and Wang, D.: The roles of climate forcing and its variability on streamflow at daily, monthly, annual, and long-term scales, *Water Resources Research*, 56, e2020WR027111, <https://doi.org/10.1029/2020WR027111>, 2020.
- Yilmaz, K. K., Gupta, H. V., and Wagener, T.: A process-based diagnostic approach to model evaluation: Application to the NWS distributed hydrologic model, *Water Resources Research*, 44, <https://doi.org/10/fpvsgb>, 2008.
- Yu, D., Yang, J., Shi, L., Zhang, Q., Huang, K., Fang, Y., and Zha, Y.: On the uncertainty of initial condition and initialization approaches in variably saturated flow modeling, *Hydrology and Earth System Sciences*, 23, 2897–2914, <https://doi.org/10.5194/hess-23-2897-2019>, 2019.
- Yu, M., Huang, Q., and Li, Z.: Deep learning for spatiotemporal forecasting in Earth system science: a review, *International Journal of Digital Earth*, 17, 2391952, <https://doi.org/10.1080/17538947.2024.2391952>, 2024.

1427 Zhang, Q., Shi, L., Holzman, M., Ye, M., Wang, Y., Carmona, F., and Zha, Y.: A dynamic
1428 data-driven method for dealing with model structural error in soil moisture data assimilation,
1429 Advances in Water Resources, 132, 103407, <https://doi.org/10.1016/j.advwatres.2019.103407>,
1430 2019.

1431 Zounemat-Kermani, M., Batelaan, O., Fadaee, M., and Hinkelmann, R.: Ensemble machine
1432 learning paradigms in hydrology: A review, Journal of Hydrology, 598, 126266,
1433 <https://doi.org/10.1016/j.jhydrol.2021.126266>, 2021.

Formatted: Header

Formatted: Indent: First line: 0"