1 **Ensembling Differentiable Process-based and Data-driven Models with**

2 **Diverse Meteorological Forcing Datasets to Advance Streamflow Simulation**

3 Peijun Li[1], Yalan Song[1], Ming Pan[2], Kathryn Lawson[1], Chaopeng Shen[1]

4 [1]Civil and Environmental Engineering, The Pennsylvania State University, University Park,

5 PA, USA

6 [2]Center for Western Weather and Water Extremes, Scripps Institution of Oceanography,

7 University of California San Diego, La Jolla, CA, USA

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

9

10 **Abstract**

11    Streamflow simulations produced by different hydrological models exhibit distinct

12 characteristics and can provide valuable information when ensembled. However, few studies

13 have focused on ensembling simulations from models with significant structural differences

14 and evaluating them under both temporal and spatial tests. Here we systematically evaluated

15 and utilized the simulations from two highly different models with great performances: a purely

16 data-driven long short-term memory (LSTM) network and a physics-informed machine

17 learning ("differentiable") HBV (Hydrologiska Byråns Vattenbalansavdelning) model (δHBV).

18 To effectively display the features of the two models, multiple forcing datasets are employed.

19 The results show that the simulations of LSTM and δHBV have distinct features and

20 complement each other well, leading to better Nash-Sutcliffe model efficiency coefficients

21 (NSE) and improved high-flow and low-flow metrics across all spatiotemporal tests, compared

22 to within-class ensembles. Ensembling models trained on a single forcing outperformed a

23 single model using fused forcings, challenging the paradigm of feeding all available data into

24 a single data-driven model. Most notably, δHBV significantly enhanced spatial interpolation

25 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 δHBV. 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.

34

**Highlights**

- Combining LSTM and δHBV with diverse forcings sets new accuracy benchmarks

- Ensembling models with one forcing outperforms merging forcings as an input

- δHBV and LSTM together always increase NSEs, especially spatial generalization

- δHBV provides valuable spatial constraints in the deterministic ensemble simulations

- δHBV and LSTM have different error characteristics that can be offset in an ensemble

41

**Keywords**

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

45

**1. Introduction**

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) 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) 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., 2025b). Conversely, process-based models, derived

59  deductively from physical laws or conceptualized views of natural systems, offer insights into

60  internal hydrological processes but may exhibit weaker performance due to structural

61  inadequacies (Li et al., 2020a, 2022; Zhang et al., 2019).

62  To combine the benefits and counteract the weaknesses of these two kinds of models,

63  many efforts have been made to incorporate physical constraints and structures into data-driven

64  models to align with fundamental physical principles, such as mass and water balances (Bandai

65  and Ghezzehei, 2021; Wang et al., 2020; Xie et al., 2021). The most seamless integration uses

66  neural networks to provide parameterizations or missing process representations for process-

67  based models (Aboelyazeed et al., 2023; Bindas et al., 2024; Feng et al., 2022; Jiang et al.,

68  2020; Kraft et al., 2022; Rahmani et al., 2023; Song et al., 2024b; Tsai et al., 2021). These

69  differentiable models (Shen et al., 2023) connect (flexible amounts of) prior physical

70  knowledge to neural networks, and have displayed many advantages, including improved

71  computational efficiency and prediction of untrained variables (Tsai et al., 2021), spatial

72  generalization (Feng et al., 2023b), and representation of extremes (Song et al., 2025b).

73  However, it is also unclear whether current differentiable models, e.g., δHBV, the

74  Hydrologiska Byråns Vattenbalansavdelning (HBV) model implemented within a

75  differentiable framework (Feng et al., 2023b; Ji et al., 2025; Shen et al., 2023; Song et al.,

76  2025b), have unique bias characteristics that are associated with the process-based parts of their

77  structures that cannot be reduced once the equations are prescribed.

78  Orthogonal to such efforts are ensemble simulations (Yu et al., 2024), which combine

79  many members with different biases and uncertainties to mitigate their respective biases in

80  deterministic predictions. Many previous studies have tried ensemble methods to improve

81  streamflow (Clark et al., 2016; Zounemat-Kermani et al., 2021) based on many factors, like

82  initial conditions (e.g., initial weights and biases in LSTM (Kratzert et al., 2018)), data used

83  for parameterization (Feng et al., 2021), and objective functions (Lin et al., 2024). These

84  studies generally use one model to generate the differences among the ensemble members.

85  Furthermore, some studies (Dion et al., 2021; Solanki et al., 2025) have utilized simulations

86  from multiple different models but are limited to process-based models, resulting in ensemble

87  simulations that are better than each individual member. Thus far, however, most studies have

88  focused on simulations from only similar models or model types, and little work has tested an

89  ensemble across the boundary of model types, particularly between data-driven, process-based,

90  and hybrid models, especially on a large number of samples. Presumably, if each model has its

91  own unique bias, data-driven and process-based models are likely to exhibit greater differences

92  due to their inherently distinct characteristics. It remains unclear whether ensembling across

93  model types should bring benefits to deterministic predictions. Furthermore, grounded in the

94  process-based model, the differentiable process-based hydrological model, such as $\delta$HBV,

95  significantly enhances performance compared to traditional process-based models, while on

96  the other hand introducing greater uncertainty regarding its potential benefits when ensembled.

97  Moreover, previous studies have primarily focused on evaluating ensemble simulations for

98  temporal predictions. However, streamflow simulation under spatial extrapolation scenarios

99  presents greater challenges, and findings from temporal tests may not be directly applicable in

100  this context.

101    It is known that the performance of any type of hydrologic model heavily depends on the

102  quality of input data, particularly meteorological forcing data (Bell and Moore, 2000; Yao et

103  al., 2020), and other inputs, like the uncertainties of initial conditions, can be mitigated via

104  warming up (Yu et al., 2019). While independent forcing datasets excel in certain aspects, they

105  each carry different error characteristics (Beck et al., 2017; Behnke et al., 2016; Newman et al.,

106  2019) and accordingly affect the hydrological models in different ways. In order to fully display

107  the different features between LSTM and $\delta$HBV, multiple forcing datasets could be considered.

108  Given the utilization of multiple forcing datasets, one could choose to use data fusion to

109    combine them into a single coherent model input (Kratzert et al., 2021; Sawadekar et al., 2025),

110    or to pass each forcing dataset through a model and then afterwards combine the multiple

111    outputs in an ensemble. It is not clear which approach is more beneficial.

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

113    questions:

114        1.  Will a cross-model-type ensemble of LSTM and $\delta$HBV improve deterministic

115            streamflow prediction more than a within-class ensemble?

116        2.  Is it better to use multiple forcings in one model or to ensemble multiple models, each

117            with a different forcing input?

118        3.  Do process-based equations bring unique value to an ensemble, especially in terms of

119            spatial generalizability?

120        The remainder of this paper is structured as follows: Sect. 2 outlines the hydrological data

121    and models used in this study, as well as the experimental design. Results and discussions are

122    presented in Sect. 3, with conclusions provided in Sect. 4.

123

124    **2. Materials and methods**

125    2.1. CAMELS hydrologic dataset

126        The Catchment Attributes and Meteorology for Large-sample Studies (CAMELS) dataset

127    (Addor et al., 2017) is widely employed for hydrological model evaluation and community

128    benchmarking. The CAMELS dataset encompasses 671 basins distributed across the

129    conterminous United States, with basin sizes ranging from 1 to 25,800 km² (median: 335 km²).

130    This standardized and publicly available dataset serves as a benchmark for evaluating various

131    hydrological models, with LSTM models trained on this dataset often serving as a reference

132    point for comparing other models (Kratzert et al., 2021). CAMELS provides basin-scale data,

133    including streamflow observations and static basin attributes, as well as forcing datasets from

134    three independent sources: Daymet (Thornton et al., 1997), North American Land Data

135    Assimilation System (NLDAS) (Xia et al., 2012), and Maurer (Maurer et al., 2002). Each of

136    the three meteorological forcing datasets operates at a daily temporal resolution, encompassing

137    precipitation, temperature, vapor pressure, and surface radiation variables, with daily

138    temperature extrema of NLDAS and Maurer supplemented from Kratzert et al. (2021). These

139    three meteorological forcing datasets have methodological distinctions in spatial resolution,

140    data generation approaches, and temporal processing (Behnke et al., 2016; Kratzert et al., 2021).

141    Exemplary plots illustrating the differences among the three meteorological forcing datasets

142    are provided in Appendix B. These features can lead to dataset-specific error characteristics

143    and make them valuable for displaying the distinct features of different model types. All model

144    inputs used in this study are detailed in Table C1.

145

146    2.2. Long short-term memory

147        As one kind of deep learning algorithm, long short-term memory (LSTM) (Hochreiter and

148    Schmidhuber, 1997) has unique structures like hidden states and gates activated by the tanh

149    and sigmoid functions (Li et al., 2023a), respectively. These features enable LSTM to excel in

150    streamflow simulation tasks (Feng et al., 2020; Kratzert et al., 2018; Nearing et al., 2024). In

151    the current benchmark framework, LSTM models are trained using dynamic atmospheric

152    forcings and static basin attributes as inputs, with streamflow as the target output, making it

153    perform well in both temporal and spatial tests (Figure 1a). In this work, for cross-group

154    comparability, we used the LSTM model and its hyperparameters as reported in Kratzert et al.

155    (2021).

156

157    2.3. Differentiable HBV model (δHBV)

158        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 (δHBV) model (Feng et al., 2023b; Ji et al., 2025; 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 δHBV 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 δHBV1.1p (Song et al., 2024b, 2025b), which is an updated version of δHBV1.0 (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~~ ~~t~~Three additional modifications are included to address high-flow simulation challenges: the use of three dynamic parameters ($\gamma$, $\beta$, $k_0$) (Song et al., 2025b); the removal of log-transform normalization for precipitation; and the adoption of the normalized squared-error loss function (Table C2) (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 δHBV are as follows:

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

$$Q = HBV(x, \theta) \tag{2}$$

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

181  where $\theta$ are the dynamic or static physical parameters, $w$ denotes the weights and biases of

182  LSTM, $x$ includes the basin-averaged meteorological forcings, such as precipitation, mean

183  temperature, and potential evapotranspiration, with $\overline{x}$ representing their normalized versions.

184  Similarly, $\overline{A_{attr}}$ consists of normalized observable basin-averaged attributes, encompassing

185  basin area, topography, climate, soil texture, land cover, and geology (Table C1). Precipitation

186  and mean temperature are from CAMELS, while potential evapotranspiration is calculated

187  ~~based on~~using the Hargreaves ~~(1994)~~(1994) method ~~using mean,~~based on maximum~~,~~ and

188  minimum temperatures along with basin latitudes, all from data described in sect. 2.1. $Q$ and

189  $Q^*$ are the streamflow simulations (model outputs) and observations (as provided in CAMELS),

190  respectively. HBV is implemented on PyTorch so it is programmatically differentiable: all

191  steps store information related to gradient calculations during backpropagation, allowing this

192  model to be trained together with neural networks in an end-to-end fashion. More details about

193  differentiable HBV can be found in previous studies ~~(Feng et al., 2022; Song et al.,~~

194  ~~2024b)~~(Feng et al., 2022; Song et al., 2024b). The details of some particularly relevant HBV

195  processes are described in Appendix A.

196

197  2.4. Experimental design

198  In this study, we trained the two models of very different types (LSTM and δHBV), each

199  with one of three meteorological forcing datasets (Daymet, NLDAS, and Maurer), resulting in

200  six corresponding streamflow simulations (Figure 1c) for each different test scenario (see sect.

201  2.5 for additional information). The training processes of LSTM and δHBV followed Kratzert

202  et al. ~~(2021)~~(2021) and Feng et al. ~~(2023b)~~(2023b), respectively. Test results and performance

203  metrics for all models are reported for the 531-basin subset that excludes those with areas larger

204  than 2,000 km² or with more than a 10% discrepancy between different basin area calculation

205 methods ~~(Newman et al., 2017)~~(Newman et al., 2017).

206     To generate ensembles, we tested various weighting strategies and ultimately employed

207 averaging to combine the six single-forcing, single-model-type simulations, as it yielded the

208 best performance. To better describe various combinations including cross-model ensembles,

209 these simulations were categorized into six groups (Table 1). A shorthand notation is used

210 throughout the remainder of this work to describe the forcing datasets and ensembles. Daymet,

211 NLDAS, and Maurer are abbreviated as superscripts 1, 2, and 3, respectively. The + symbol is

212 used to group model types being ensembled, while superscript clustering (e.g., $^{12}$ or $^{123}$) is used

213 to group the meteorological forcing types being ensembled, with parentheses indicating that

214 the superscripts apply to all model types within. For example, $(LSTM + \delta HBV)^{123}$ could be

215 explicitly written as $LSTM^1 + LSTM^2 + LSTM^3 + \delta HBV^1 + \delta HBV^2 + \delta HBV^3$. To compare

216 two different strategies to utilize the multiple meteorological forcing datasets and to benchmark

217 against the previously highest performance, we additionally trained a single LSTM model using

218 all three forcing datasets as simultaneous inputs as done by Kratzert et al. ~~(2021)~~(2021),
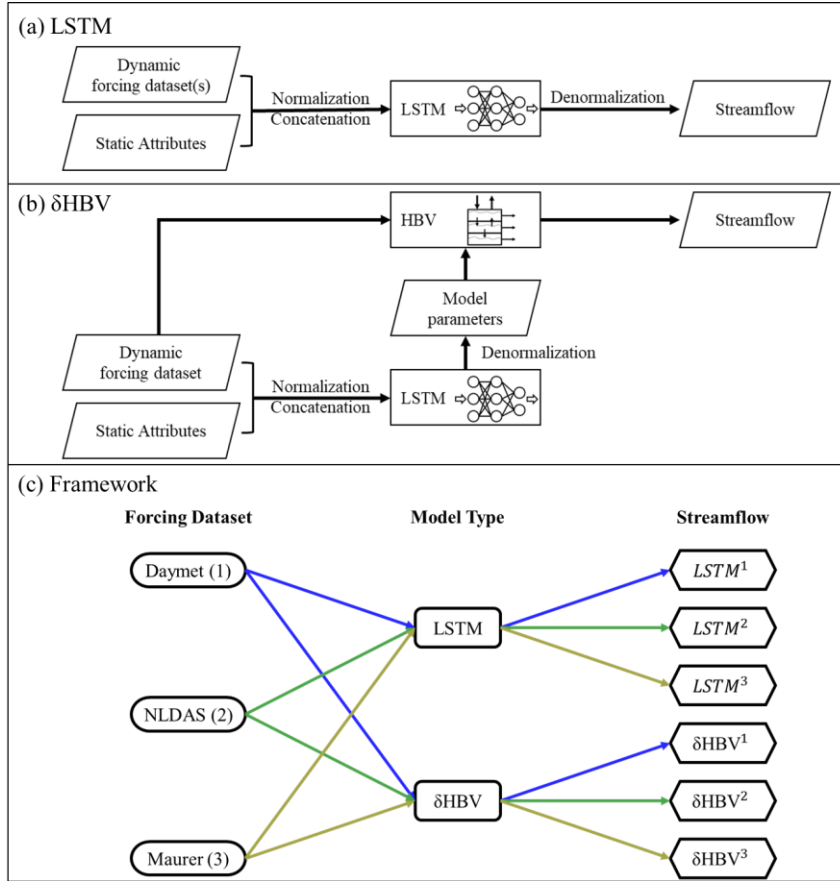
219 referred to as $LSTM^{multi}$ (the last row in Table 1).

220

*Figure 1. (a) The LSTM structure, (b) the δHBV 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),*
*respectively.*

225 *Table 1. (a) The six groups of streamflow simulations, and (b) the streamflow simulation via*
226 *LSTM based on a different strategy, in which three meteorological forcing datasets were*
227 *combined as a single set of inputs ~~(Kratzert et al., 2021)~~(Kratzert et al., 2021). Superscripts 1,*
228 *2, and 3 denote Daymet, NLDAS, and Maurer, respectively. The ensemble across forcings*
229 *("ef") superscript indicates an ensemble of model simulations, each of which uses a different*
230 *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 | $LSTM^{1}$, $LSTM^{2}$, $LSTM^{3}$ |
| $\delta$HBV | $\delta HBV^{1}$, $\delta HBV^{2}$, $\delta HBV^{3}$ |
| LSTM+$\delta$HBV | $(LSTM+\delta HBV)^{1}$, $(LSTM+\delta HBV)^{2}$, $(LSTM+\delta HBV)^{3}$ |
| $LSTM^{ef}$ | $LSTM^{12}$, $LSTM^{13}$, $LSTM^{23}$, $LSTM^{123}$ |
| $\delta HBV^{ef}$ | $\delta HBV^{12}$, $\delta HBV^{13}$, $\delta HBV^{23}$, $\delta HBV^{123}$, |
| $(LSTM+\delta HBV)^{ef}$ | $(LSTM+\delta HBV)^{12}$, $(LSTM+\delta HBV)^{13}$, $(LSTM+\delta HBV)^{23}$, $(LSTM+\delta HBV)^{123}$ |

| (b) Using forcing datasets as simultaneous inputs to an LSTM | | |
|---|---|---|
| **Streamflow Simulation** | **Model Type** | **Meteorological Forcing Dataset** |
| $LSTM^{multi}$ | LSTM | Daymet, NLDAS, Maurer |

231
232

233 2.5. Evaluation scenarios and criteria

234 The above cases were comprehensively evaluated for performance in temporal

235 extrapolation (Feng et al., 2022; Kratzert et al., 2018)(Feng et al., 2022; Kratzert et al., 2018),

236 as well as two types of spatial generalization: prediction in ungauged basins (PUB) (Feng et

237 al., 2023b; Kratzert et al., 2019)(Feng et al., 2023b; Kratzert et al., 2019), and prediction in

238 ungauged regions (PUR) (Feng et al., 2021, 2023b):

239 • **Temporal Test**: Models were trained using data from all basins and tested across

240    different periods.

241 • **PUB Test**: Models were trained on randomly selected subsets from all basins and

242    tested on the remaining basins during the same time period.

243 • **PUR Test**: Different from the PUB test, basins were grouped into continuous regions,

244    one of which was selected to comprise the group of testing basins while the others

245    were used for training.

246 Temporal generalization is generally considered to be the easiest of these tests. In terms

247 of spatial generalization, which approximates data-sparse scenarios, the PUB test is an example

248 of spatial interpolation, whereas the PUR test involves spatial extrapolation. The PUR test is

249 widely regarded as the most challenging and may therefore produce findings that differ

250 significantly from those in other scenarios. In this study, all basins were divided into 10

251 spatially stratified groups for the PUB test and 7 fully disjoint regional groups for the PUR test

252 (Table 2) in the same way as Feng et al. (2023b)(2023b). The spatial extent of the 7 regions for

253 the PUR test is also shown in Figure 3(c1-c2). Therefore, we conducted 10 rounds for the PUB

254 test and 7 rounds for the PUR test, with a different group held out for testing in each round.

255 Model performance was evaluated after concatenating the test results for all basins.

256

257 *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 |

258 [a]*δHBV training followed Feng et al.* ~~(2023b)~~*(2023b) using all 671 CAMELS basins, while*
259 *LSTM training followed Kratzert et al* ~~(2021)~~*(2021) using the selected 531-basin subset. Test*
260 *results and performance metrics for all models are reported for the 531 basins.*
261 [b]*Each hydrological year spans from October 1st to September 30th of the following year.*
262 [c]*In the PUB and PUR tests, models are run for 10 and 7 rounds, respectively, with the group*
263 *held out for testing changed in each round. The simulation performance was evaluated after*
264 *concatenating the test results for all basins.*
265

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

14

279  5 and 21 hours, respectively, using a single NVIDIA Tesla V100 GPU.

280  We calculated several well-established performance metrics: Nash-Sutcliffe model

281  efficiency coefficient (*NSE*) ~~(Nash and Sutcliffe, 1970)~~(Nash and Sutcliffe, 1970), Kling-

282  Gupta model efficiency coefficient (*KGE*) ~~(Kling et al., 2012)~~(Kling et al., 2012), percent bias

283  (*PBIAS*), and root-mean-square error (*RMSE*). We also considered *RMSE* values for high (top

284  2% "peak" flow, *highRMSE*), low (bottom 30% "low" flow, *lowRMSE*), and mid-range (the

285  remaining flow, *midRMSE*) flow conditions ~~(Yilmaz et al., 2008)~~(Yilmaz et al., 2008). These

286  metrics were computed for each basin and aggregated into error bars and cumulative density

287  functions (CDFs). For brevity, the main text primarily reports NSE values, and other metric

288  values are provided in Appendices D and E. Furthermore, we use the spread values ~~(Li et al.,~~

289  ~~2021; Reichle and Koster, 2003)~~(Li et al., 2021; Reichle and Koster, 2003) to investigate

290  ensemble variability and explore model complementarity. Detailed descriptions of these

291  metrics and their calculations are available in Table C2.

292

293  **3. Results and discussion**

294  3.1. Temporal extrapolation

295  For the temporal test, in which models were trained and tested on the same basins but in

296  different time periods, we found that cross-model-type ensembles noticeably surpassed the

297  within-class ensembles when other conditions were the same, with small uncertainties as shown

298  by the error bars in Figure 2. With a single forcing dataset, the median NSE was elevated from

299  ~0.735 for LSTM to ~0.79 with δHBV added, though δHBV performance was similar to LSTM

300  (~0.74 under Daymet). Even after LSTM achieved very high performance when its simulations,

301  each derived separately from different meteorological forcing datasets, were ensembled (ef =

302  123, ~0.808), adding δHBV still improved the results to ~0.818. This finding was robust for

303  all different combinations of the tested meteorological forcing datasets. Conversely, adding

304  LSTM also helped to improve δHBV ensembles. These results highlight the benefits of the

305  cross-model-type ensemble framework and indicate distinct simulation features for each model

306  type. LSTM is a data-driven method that has low bias and large variance. ~~Data errors (Li et al.,~~

307  ~~2020b)~~Data errors (Li et al., 2020b), different sampling strategies ~~(Nai et al., 2024)~~(Nai et al.,

308  2024), or even different weight initializations ~~(Narkhede et al., 2022)~~(Narkhede et al., 2022)

309  can lead to substantively different outcomes. Conversely, δHBV may have a smaller variance

310  but a larger bias due to the fixed HBV formulation ~~(Moges et al., 2016)~~(Moges et al., 2016) for

311  some scenarios like low flows ~~(Feng et al., 2023b; Song et al., 2024b)~~(Feng et al., 2023b; Song

312  et al., 2024b) or in basins with significant water uses ~~(Song et al., 2024a)~~(Song et al., 2024a).

313  These errors with varying characteristics from different model classes can partially offset each

314  other in an ensemble. On a side note, δHBV models seem more reliant on the quality of the

315  forcing data, as shown in Figure 2. δHBV with the Maurer and NLDAS forcing datasets

316  generally performs worse than it does with Daymet, which has lower biases. However, even in

317  those cases, the combination of LSTM and δHBV was still better than LSTM alone, attesting

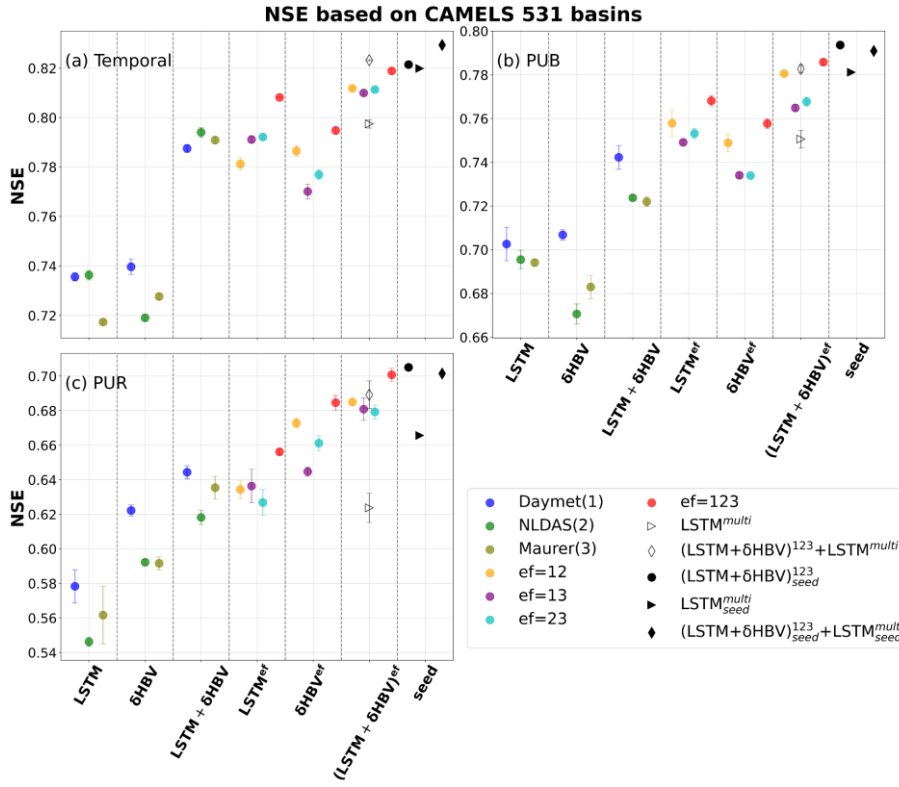318  to the robustness of these benefits.

319

320

*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.*
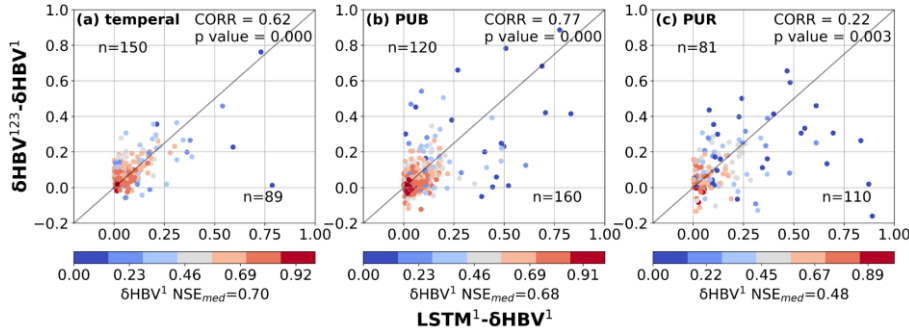
*Figure 3. Scatter plots comparing the performance differences between hydrological models for the basins where LSTM outperformed δHBV (the basins where δHBV outperformed are not shown in this plot). The x-axis represents the NSE differences between $LSTM^1$ and $δHBV^1$ ($LSTM^1$ - $δHBV^1$), while the y-axis shows the NSE differences between $δHBV^{123}$ and $δHBV^1$ ($δHBV^{123}$- $δHBV^1$). 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 generally not be subtracted. Here the purpose is only to confirm that basins where LSTM outperforms δHBV 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 $δHBV^2$ 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

18

359 same basins, this also explains how LSTM[1] surpasses δHBV[1] in some basins with poorer-

360 quality forcings. In contrast to δHBV, LSTM has the innate ability to shift information in time

361 and moderately adjust the input scale. Moving from temporal validation to PUB to PUR

362 scenarios, the advantages of diverse forcing datasets appear to diminish, as evidenced by the

363 decreasing ratio of points above versus below the diagonal line, since the forcing error patterns

364 remembered by LSTM may not generalize well in space (discussed in more detail in sect. 3.2).

365

366     Ensembling streamflow simulations from different meteorological forcing datasets

367 demonstrates certain advantages over the previous approach of simultaneously sending

368 multiple forcings into a data-driven model like LSTM ~~(Kratzert et al., 2021)~~(Kratzert et al.,

369 2021). Ensembling LSTM simulations each using a single forcing dataset ($LSTM^{123}$) resulted

370 in an NSE value of 0.8082, higher than that of 0.7974 from feeding multiple forcing datasets

371 into a single LSTM ($LSTM^{multi}$). This difference was more pronounced in the cross-model-

372 type ensemble, after including δHBV, compared to the previous within-class ensemble, and

373 particularly notable for the spatial generalization tests (to be discussed in more detail in Sect.

374 3.2~~), with~~). The corresponding specific ~~metric values~~performance metrics are summarized in

375 Tables D1–D5, with seasonal evaluations provided in ~~Tables D1–D5~~Figure D3. These results

376 indicate that the trained LSTM in $LSTM^{multi}$ may be overfit~~ted~~ to the significant redundant

377 information in these three forcing datasets, and that ~~only~~ LSTM models alone cannot fully

378 exploit the information hidden in the multiple forcing datasets. Training separate ensemble

379 members via different nonlinear hydrological processes, on the other hand, seems to allow

380 different bias features to emerge with separate forcing datasets, accordingly mitigating them

381 during the subsequent ensembling process.
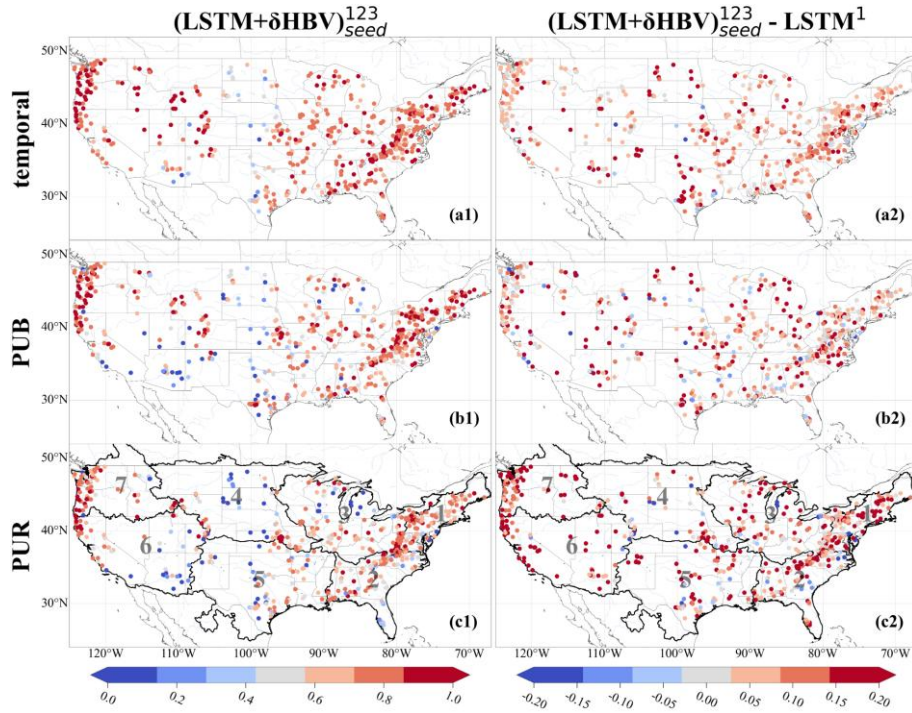
382

19

*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, $(LSTM + \delta HBV)^{123}_{seed} + LSTM^{multi}_{seed}$, 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, $LSTM^{multi}_{seed}$, 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~~improved, the gaps between them did not change much (Figure 5 ~~and~~; Table D5 for 10 seeds, Table D4 for 3 seeds). In particular, the gap between $(LSTM + \delta HBV)^{123}_{seed} + LSTM^{multi}_{seed}$ and $(LSTM + \delta HBV)^{123}_{seed}$ or $LSTM^{multi}_{seed}$ 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 $(LSTM + \delta HBV)^{123}$ with additional random seeds. It was noteworthy that while the $(LSTM + \delta HBV)^{123}$ 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 $LSTM^{multi}_{seed}$ variant into $(LSTM + \delta HBV)^{123}_{seed} + LSTM^{multi}_{seed}$, 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.
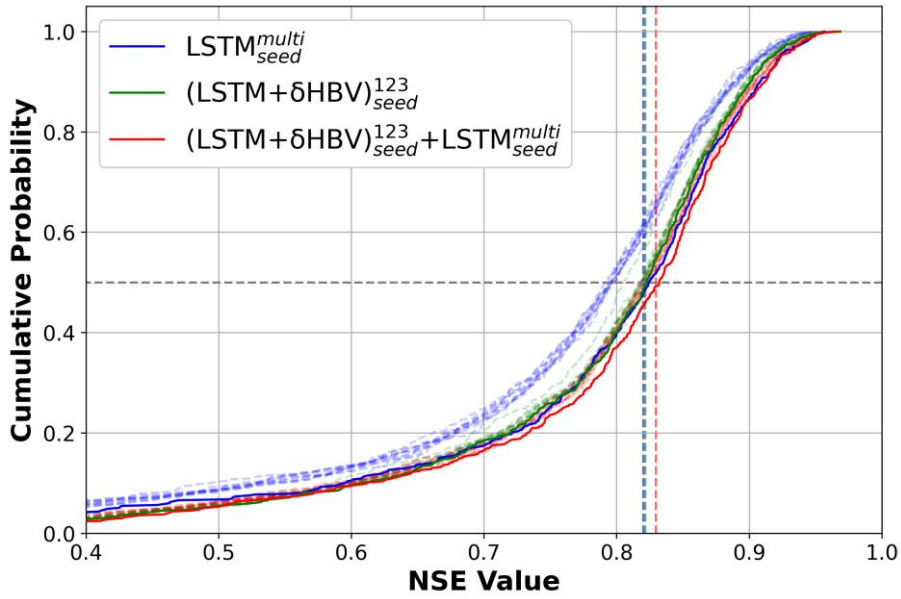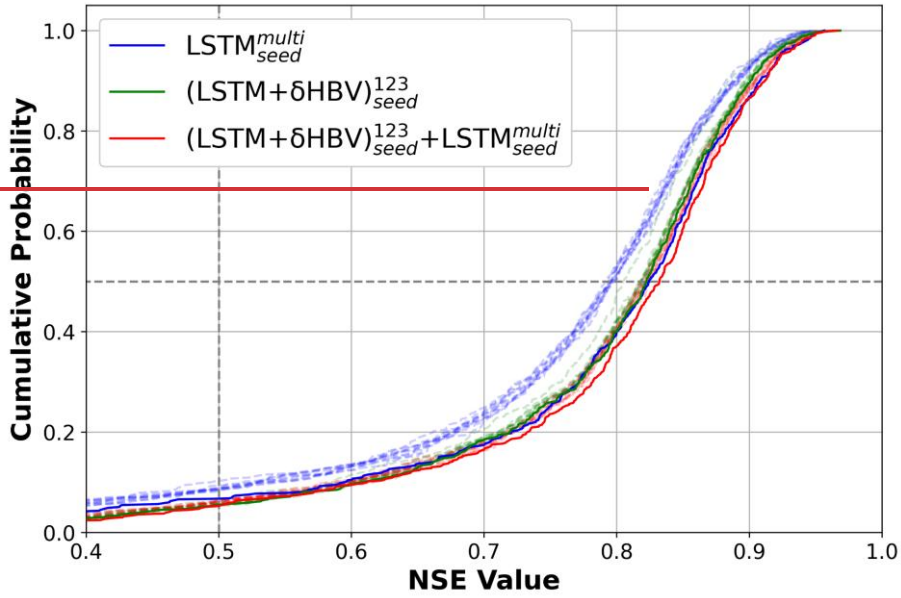
419



420

*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. The median NSE values computed with 3 random seeds are also indicated by vertical dashed and translucent lines in the corresponding colors.*

3.2. Spatial generalization

It is clear that cross-model-type ensembling and the incorporation of δHBV 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, δHBV-only individual simulations or ensembles consistently outperformed LSTM-only counterparts in the PUR test. Furthermore, adding δHBV to the same-model-type LSTM ensembles improved median NSE by 0.02-0.03 for PUB. The role of δHBV 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 δHBV is also illustrated by the optimized weights shown in Figure E1, which were estimated using a genetic algorithm with streamflow observations from the test periods. These weights are presented solely to 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 $LSTM^{multi}$ regardless of whether ~~or not we employed~~ multiple random seeds were employed to form an ensemble. As such, we can conclude that the inclusion of a differentiable process-based model like δHBV 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

23

450 (Figure 4-b2), and PUR (Figure 4-c2) test. The most significant improvements due to the

451 ensemble were concentrated on the center of the Great Plains along with the midwestern US,

452 while the eastern US was moderately improved, suggesting data uncertainty is a larger issue in

453 the central and midwestern US. The Great Plains have historically had poor performance for

454 all kinds of models (Mai et al., 2022)(Mai et al., 2022) and even the ensemble model had NSE

455 values of only 0.3-0.4 for many of the basins there, although this still marked significant

456 improvements over LSTM[1] (Figure 4-a2, -b2, -c2). Some western basin NSE values were

457 elevated by more than 0.15 for the temporal test (Figure 4-a2) and even more for PUB and

458 PUR. Meteorological stations are generally sparse on the Great Plains, and an ensemble seems

459 to be an effective way to leverage the different forcing datasets that are available. The poor

460 performances in some basins highlight some remaining deficiencies in current models, which

461 clearly cannot fully consider the heterogeneities of different basins; thus, multiscale

462 formulations that resolve such heterogeneities may have advantages (Song et al., 2024a)(Song

463 et al., 2024a).

464     To investigate why ensembles outperformed single-model, single-forcing approaches, we

465 compared their temporal, PUB, and PUR test simulation time series against observations for

466 531 basins (Figure 6). Analysis of averaged hydrological year data revealed that while

467 individual ensemble members using single-source forcing datasets performed similarly for

468 easily simulated periods, they showed significant divergence during challenging periods,

469 particularly peak flows. This divergence stems from distinct systematic errors inherent to

470 different model types and forcing datasets. Notably, LSTM-based simulations alone proved

471 insufficient in generating adequate spread to capture these divergent points. By averaging

472 individual model outputs and stabilizing uncertainties, ensemble simulations achieved effective

473 and robust performance across all conditions, which can be shown via the metric highRMSE

474 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 ensure that each member develops distinct 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.
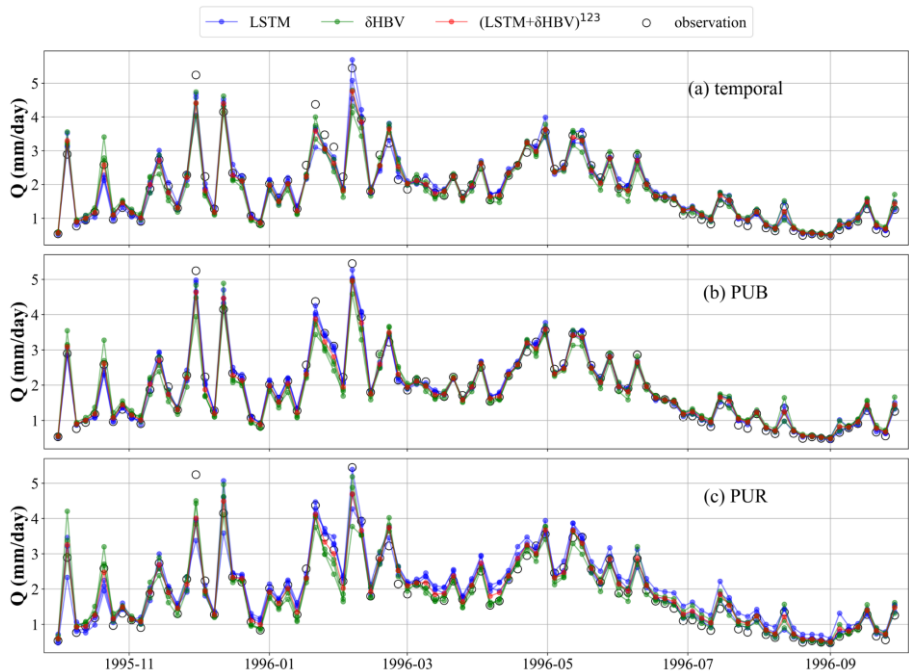


*Figure 6. Comparisons between multi-basin-averaged streamflow observations and simulations across 531 basins. The time series points are displayed at four-day intervals for clarity and conciseness. Ensemble members based on the same model (LSTM or δHBV) but driven by different forcing datasets are shown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

3.3 Ensemble variability and robustness analysis

　　Although δHBV (median spread 0.61) exhibits lower spreads than LSTM (mean spread

491    0.72), their combination increases the ensemble spreads, thereby enhancing diversity (Figure

492    7). This pattern holds across the temporal, PUB, and PUR tests. Ensemble effectiveness

493    depends on the diversity of model behaviors and their distinct error characteristics.

494    Consequently, larger spreads ~~contribute to~~are generally associated with greater ensemble

495    benefits. Figure ~~D3~~D4 further demonstrates that δHBV+LSTM exhibits larger spreads than
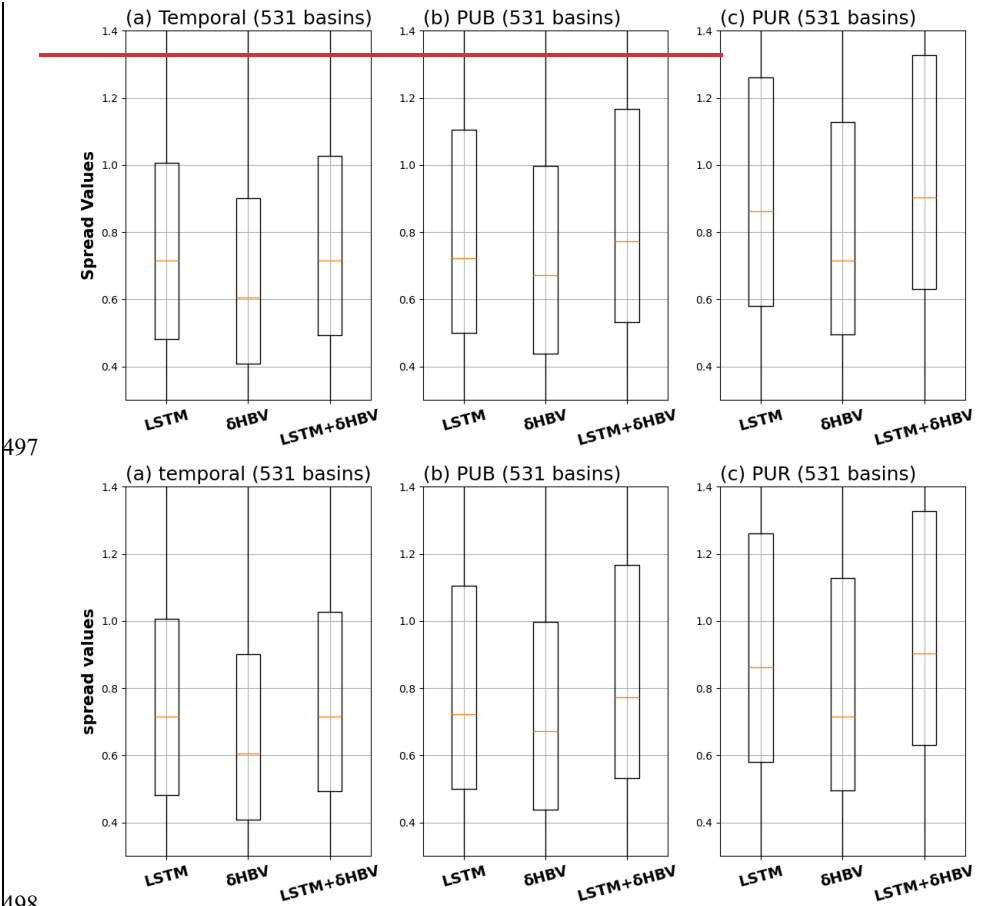
496    LSTM in most basins.

497

498



499    *Figure 7. Spread values (Table C2) of each model for LSTM, δHBV, and LSTM + δHBV due*

500    *to different meteorological forcings and random seeds across temporal, PUB, and PUR tests.*

501

As the warming signal is already clear across most basins under any forcing across the periods of simulation (Figure ~~D4~~D5), 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)~~(Feng et al., 2023b). Introducing a 10% precipitation perturbation (multiplying precipitation by 1.1) slightly reduced performance for both models as expected (Figures D~~5~~6a and D~~5~~6b), but ensemble benefits remained robust across models despite the perturbation.

Training sample size, dynamic parameter choices, and lookback windows ~~all have~~exert only a limited impact on our conclusions. δHBV shows limited sensitivity to sample size, with similar results when trained on 531 versus 671 basins (Figure D~~5~~6c). Regarding parameter uncertainties, fixing one δHBV parameter ($k_0$) as static increased structural errors and reduced performance (Figure D~~5~~6d), 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 D~~5~~6e). 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. 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

27

include using the simulations from these 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 δHBV models exhibit limitations in regions with significant anthropogenic impacts, such as ~~like~~ dam presence, as well as arid climatic and highly~~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) or integrating specialized modules, such as reservoirs (Hanazaki et al., 2022; West et al., 2025). Compared with LSTM, δHBV is more sensitive to precipitation biases. For example, the differences between δHBV simulations under different forcing datasets were generally larger than those for LSTM, and δHBV using the Daymet forcing dataset showed largely better performance than with the other two forcing datasets, which indicates that δHBV may not be able to fit different forcing datasets well. Therefore, many potential structural optimizations can be implemented to improve δHBV. Our analysis provided corroborating evidence that forcing error is an important reason why LSTM can outperform δHBV 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. 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 $\delta HBV^{multi}$ 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 $\delta HBV^{multi}$ 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 $\delta HBV$ has lower variance and potentially higher bias. As a result, random seeds would likely not create large enough perturbations for $\delta HBV$ and wouldn't bring the benefits seen with $LSTM_{seed}^{multi}$. To achieve an equivalent perturbation level for $\delta HBV$, 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)~~(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 conclusions

This study comprehensively analyzes ensemble combinations of two advanced model

578  types (LSTM and δHBV), each with distinct mechanisms, for streamflow simulation across

579  531 basins in the US. Three meteorological forcing datasets (Daymet, NLDAS, and Maurer)

580  are employed to fully capture the characteristics of the two models. Their applications are also

581  tested in two distinct ways: (1) by feeding all diverse forcing datasets simultaneously into a

582  single LSTM model, and (2) by ensembling the outputs of multiple LSTM models, each trained

583  separately using a single forcing dataset. The performance of ensemble simulations was

584  evaluated under three distinct testing scenarios (temporal, PUB, and PUR tests), surpassing the

585  previous highest performances. Our findings enhance the understanding of how to effectively

586  utilize diverse model types and multi-source datasets to improve streamflow simulations. The

587  principal conclusions are:

588  (1) Cross-model-type ensembles (LSTM+δHBV) consistently outperformed single-
589  model approaches across all test scenarios, setting new performance benchmarks on
590  the CAMELS dataset. These ensembles demonstrated the complementarity of data-
591  driven (LSTM) and physics-informed (δHBV) approaches in capturing diverse
592  hydrological behaviors.

593  (2) Ensembling models trained on different forcing datasets proved more effective than
594  using multiple forcing datasets as simultaneous inputs to a single model. This suggests
595  that separate training allows each model to capture unique features contained in each
596  forcing dataset, which can then be effectively leveraged in the ensemble.

597  (3) δHBV provided significant benefits to ensemble simulations on spatial generalization.
598  Ensembling LSTM with δHBV showed increasing benefits as generalization
599  challenges increased, from temporal to spatial interpolation (PUB) to spatial
600  extrapolation (PUR) tests. This underscores the value of physics-informed constraints
601  in improving model transferability to ungauged basins and regions.

602  (4) While ensemble methods significantly improved overall performance, they did not

603     fully mitigate consistent deficiencies in certain challenging areas (e.g., regions with

604     high dam density or heterogeneous hydrogeological conditions). This indicates areas

605     for future model development.

606     These findings have important implications for hydrological modeling and water

607 resources management. The improved accuracy and spatial generalization of our ensemble

608 approach can enhance streamflow predictions, benefiting water resources planning and

609 management, particularly in data-scarce regions. Our results also suggest that future

610 hydrological model development should focus on combining data-driven and physics-based

611 approaches to improve model generalizability across diverse conditions. The superior

612 performance of ensembling models with different forcing datasets over using merged forcings

613 as a single input highlights the risk of indiscriminately feeding all available data into one data-

614 driven model. While computational demands certainly require consideration, the potential

615 improvements in prediction accuracy offer significant value for both research and operational

616 applications. Future work should focus on refining these ensemble techniques, addressing

617 model limitations in challenging regions, and exploring ensemble implementation in

618 operational settings.

**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).~~

The Hydrologiska Byråns Vattenbalansavdelning (HBV) model (Aghakouchak and Habib, 2010; Beck et al., 2020; Bergström, 1976, 1992) is a simple yet effective bucket-type hydrologic model that simulates hydrologic components including snow water equivalent, soil moisture, groundwater storage, evapotranspiration, quick flow, baseflow, and total streamflow. In the following, we describe these processes in detail with their corresponding equations. Uppercase letters denote state variables, while lowercase letters denote parameters. The overall water balance is expressed as Equation (S1).

$$EP - AE - Q_t = SN + SM + \text{~~UR + LR~~}SUZ + SLZ + LAKE \tag{S1}$$

where $EP$ is effective precipitation, $AE$ is ~~the~~ actual evapotranspiration, $Q_t$ is ~~the~~ total simulated runoff, $SN$ is snow storage, $SM$ is soil ~~water~~moisture storage, ~~$UR$ is~~$SUZ$ and $SLZ$ are the upper ~~reservoir water level, $LR$ is the~~ and lower ~~reservoir water level~~groundwater storages, respectively, and $LAKE$ ~~is the~~represents lake ~~level~~storage (omitted in this study).

First, effective precipitation ($EP$) is ~~separa~~partition~~t~~ed into ~~liquid~~rain ($RN$) and ~~solid~~snow ($SN$) components based on the air temperature ($T$) relative to ~~the~~a threshold temperature ($tt$)~~ as~~:

$$RN = EP \ if \ T \geq tt \tag{S2}$$

$$SN = EP \ if \ T < tt \tag{S3}$$

Snow ($SN$) accumulates in the snowpack ($SNP$), while ~~the~~ snowmelt ($SNM$) happens when $T \geq tt$, which is calculated ~~using~~based on a ~~temperature-dependent~~ melt ~~rate~~factor ($cfm$~~.~~) and the temperature difference ($T - tt$). The computed snowmelt ($SNM$) is ~~limited to~~constrained by the available snowpack ($SNP$~~, and any excess melt contributes to meltwater ($MW$) as~~).

$$\text{~~}SNP = SNP + SN\text{~~} \tag{S4}$$

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

$$\text{~~}MW = MW + SNM\text{~~} \tag{S6}$$

$$\text{~~}SNP = SNP - SNM\text{~~}SNM$$
$$= min[max(cfm \cdot (T - tt), 0), SNP] \tag{~~S7~~S4}$$

644 ~~Some of this~~ The snowmelt (*SNM*) contributes to meltwater (*MW*~~) refreezes based on a~~

645 ~~refreezing parameter (*cfr*) and the temperature difference from the threshold, returning to~~ ),

646 while the snowpack (*SNP*~~). The amount of refrozen water~~) is ~~labeled~~updated as ~~*FRZ*~~:

$$MW = MW + SNM\text{~~}RFZ\text{~~}$$

$$= \begin{cases} MW & cfr \cdot cfm \cdot (tt - T) \geq MW \\ cfr \cdot cfm \cdot (tt - T) & T < tt, cfr \cdot cfm \cdot (tt - T) < MW \\ 0 & T \geq tt \end{cases} \tag{~~S8~~S5}$$

$$SNP = SNP + \text{~~}RFZ\text{~~}SN - SNM \tag{~~S9~~S6}$$

$$\text{~~}MW = MW - RFZ\text{~~} \tag{~~S10~~}$$

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

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

649 A portion of the meltwater (*MW*) may refreeze when *T* < *tt*, controlled by the refreezing

650 parameter (*cfr*):

$$\text{~~}IF = \begin{cases} MW - cwh * SNP & MW - cwh * SNP \geq 0 \\ 0 & MW - cwh * SNP < 0 \end{cases}RFZ\text{~~}$$
$$= min[max(cfr \cdot cfm \cdot (tt - T), 0), MW] \tag{~~S11~~S7}$$

$$SNP = SNP + RFZ \tag{S8}$$

$$MW = MW - IF\text{~~}Z\text{~~} \tag{~~S12~~S9}$$

651 ~~The fraction of soil moisture relative to the field capacity (*fc*) determines the soil wetness,~~

652 ~~which modulates the amount of water recharged into the soil (*SP*). Then soil moisture (*SM*) is~~

653 ~~updated based on the infiltration of meltwater (*IF*), rain (*RN*), and the amount of recharged~~

654 ~~water (*SP*) as~~

655 The remaining meltwater (*MW*) exceeding the snowpack's liquid water holding capacity (*cwh* ·

656 *SNP*) infiltrates into the soil (*IF*), with the remainder retained in *MW*:

$$IF = max(MW - cwh \cdot SNP, 0)\text{~~}SP = \left(\frac{SM}{fc}\right)^\beta \cdot (IF + RN)\text{~~} \tag{S1~~3~~0}$$

$$\text{~~}SM = SM +\text{~~}MW = MW - IF\text{~~} + RN - SP\text{~~} \tag{S1~~4~~1}$$

657 The ~~excess water, above~~fraction of soil moisture (*SM*) relative to the field capacity (~~*IF*<sub>air</sub>), is~~

658 ~~calculated and subsequently removed from~~*fc*), raised to the ~~soil moisture storage as~~power index

659 *β*, modulates shallow seepage (*SP*) according to the available water (*IF* + *RN*):

$$SP = \left(\frac{SM}{fc}\right)^\beta \text{~~}IF_{air} = \begin{cases} SM - fc & if\ SM \geq fc \\ 0 & if\ SM < fc \end{cases}\text{~~}(IF + RN) \tag{S1~~5~~2}$$

$$SM = SM \; \cancel{- IF_{dir}} + IF + RN - SP \tag{S1$\cancel{6}$3}$$

660 ~~Actual evapotranspiration (*AE*) is determined by an evaporation factor (*PEC*), which depends~~
661 ~~on the soil moisture, a shape parameter (*λ*), a parameter (*lp*), and field capacity (*fc*) for~~
662 ~~evapotranspiration. This factor limits the actual evapotranspiration (*AE*) to both the potential~~
663 ~~evapotranspiration (*PE*) and the available soil moisture.~~
664 Excess soil water above the field capacity contributes to direct infiltration ($IF_{dir}$):

$$\cancel{PEC = \begin{cases} \left(\dfrac{SM}{lp\cdot fc}\right)^{\lambda} & if\ 0 \le \left(\dfrac{SM}{lp\cdot fc}\right)^{\lambda} < 1 \\ 0 & if\ S\left(\dfrac{SM}{lp\cdot fc}\right)^{\lambda} < 0 \\ 1 & if\ S\left(\dfrac{SM}{lp\cdot fc}\right)^{\lambda} \ge 1 \end{cases} IF_{dir}} \tag{S1$\cancel{7}$4}$$

$$= \max(SM - fc, 0)$$

$$\cancel{AE = \begin{cases} PE\cdot PEC & if\ SM \ge PE\cdot PEC \\ SM & if\ SM < PE\cdot PEC \end{cases}} \tag{S18}$$

$$SM = SM - \cancel{AE}\,IF_{dir} \tag{S1$\cancel{9}$5}$$

665 ~~Capillary rise (*CP*) from the lower soil zone (*SLZ*) is governed by a parameter (*c*), which~~
666 ~~determines the amount of water moving upward based on the soil moisture content. This~~
667 ~~capillary flow replenishes the soil moisture, while groundwater interactions occur through~~
668 ~~recharge processes in the upper (*SUZ*) and lower (*SLZ*) groundwater zones.~~
669 Actual evapotranspiration (*AE*) is estimated as the product of potential evapotranspiration (*PE*)
670 and an evapotranspiration coefficient (*PEC*). The *PEC* depends on soil moisture storage (*SM*),
671 field capacity (*fc*), a shape parameter (*λ*), and a threshold parameter (*lp*).

$$\cancel{CP = \begin{cases} SLZ & if\ SLZ < c\cdot SLZ\cdot\left(1-\dfrac{SM}{fc}\right) \\ c\cdot SLZ\cdot\left(1-\dfrac{SM}{fc}\right) & if\ SLZ \ge c\cdot SLZ\cdot\left(1-\dfrac{SM}{fc}\right) \end{cases} PEC} \tag{S$\cancel{20}$S16}$$

$$= \min\left[1, \max\left(0, \left(\frac{SM}{lp\cdot fc}\right)^{\lambda}\right)\right]$$

$$AE = \min(PE\cdot PEC, SM) \qquad SM = SM + CP \tag{S$\cancel{2}$17}$$

$$\cancel{SLZ = \begin{cases} SLZ - CP & if\ SLZ \ge CP \\ 0 & if\ SLZ < CP \end{cases}} SM = SM - AE \tag{S$\cancel{22}$S18}$$

672 ~~Excess recharge (*SP* and *IF$_{dir}$*) from the soil enters the upper zone, where it either percolates~~
673 ~~to the lower zone (*PERC*) based on a constant rate (*prc*) or contributes to direct runoff (*Q$_0$*)~~
674 ~~when it exceeds the upper zone threshold (*uzl*). The generated flow is modeled using~~

parameters ($k_0$, $k_1$, $k_2$) governing flow from the upper and lower zones. Each of these flows contributes to runoff ($Q_0$Capillary rise ($CP$) from the lower zone ($SLZ$) replenishes $SM$, controlled by a coefficient (c) and constrained by the soil moisture deficit:

, $Q_1$, $Q_2$), and their respective contributions to streamflow ($Q_t$) are modeled over time.

$$CP = \min\left[c \cdot SLZ \cdot \left(1 - \frac{SM}{fc}\right), SLZ\right] SUZ = SUZ + SP + IF_{dir} \qquad \text{(S23S19)}$$

$$SM = SM + CP PERC = \begin{cases} prc & if\ SUZ \geq prc \\ SUZ & if\ SUZ < prc \end{cases} \qquad \text{(S240)}$$

$$SUZ = SUZ - PERC \qquad \text{(S25)}$$

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

$$SUZ = SUZ - Q_0 \qquad \text{(S27)}$$

$$Q_1 = SUZ \cdot k_1 \qquad \text{(S28)}$$

$$SUZ = SUZ - Q_1 \qquad \text{(S29)}$$

$$SLZ = SLZ + PERC \qquad \text{(S30)}$$

$$Q_2 = SLZ \cdot k_2 \qquad \text{(S31)}$$

$$SLZ = SLZ - Q_2 CP \qquad \text{(S321)}$$

$$Q_t = Q_0 + Q_1 + Q_2 \qquad \text{(S33)}$$

Recharge from the soil, consisting of shallow seepage ($SP$) and direct infiltration ($IF_{dir}$), enters the upper groundwater zone ($SUZ$). Water in the upper zone either percolates to the lower groundwater zone ($SLZ$) at a constant percolation rate ($prc$) or contributes to direct runoff ($Q_0$) when the upper zone (SUZ) exceeds a threshold ($uzl$). Flow from the upper and lower zones is computed using linear reservoir formulations, with parameters $k_0$, $k_1$, $k_2$ controlling the respective runoff components $Q_0$, $Q_1$, $Q_2$. The total simulated streamflow ($Q_t$) is then computed as the sum of these components.

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

$$SUZ = SUZ + SP + IF_{dir} Q_t^* = \int_0^{tmax} \xi(s: \theta_\alpha, \theta_t) \cdot Q(t - s) ds \qquad \text{(S34S22)}$$

$$\cancel{\xi(s:\theta_\alpha,\theta_\tau) = \frac{1}{\Gamma(\theta_\alpha)\theta_\tau^{\theta_\alpha}} t^{\theta_\alpha-1} e^{-\frac{t}{\theta_\tau}}} \quad PERC = \min(prc, SUZ) \qquad \text{(S23)}$$

$$SUZ = SUZ - PERC \qquad \text{(S24)}$$

$$Q_0 = \max[k_0 \cdot (SUZ - uzl), 0] \qquad \text{(S25)}$$

$$SUZ = SUZ - Q_0 \qquad \text{(S26)}$$

$$Q_1 = SUZ \cdot k_1 \qquad \text{(S27)}$$

$$SUZ = SUZ - Q_1 \qquad \text{(S28)}$$

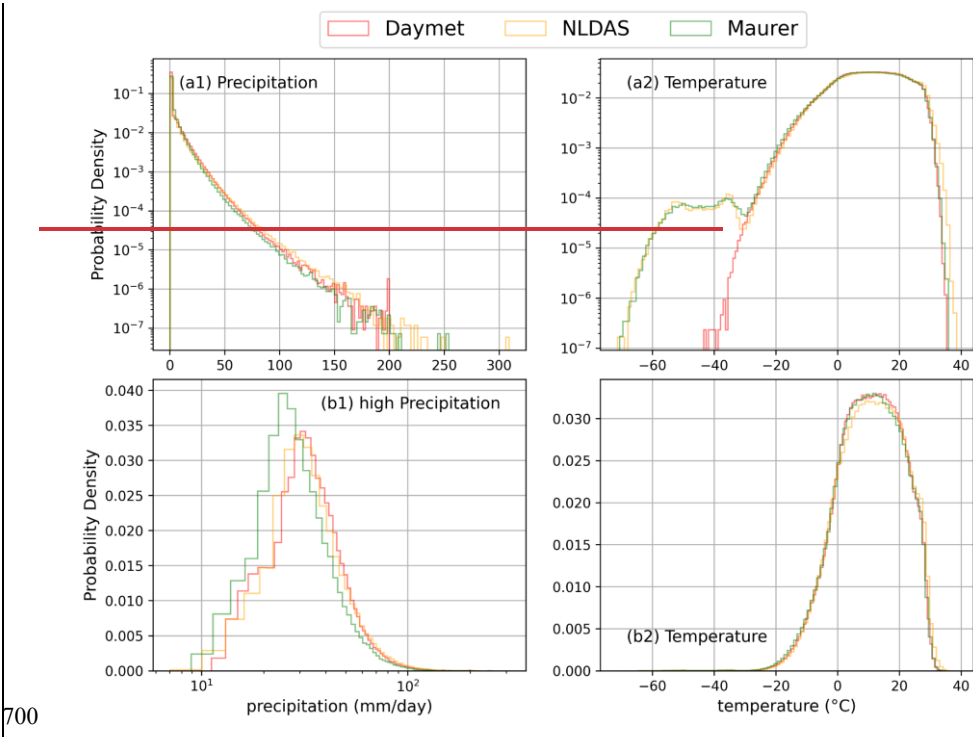$$SLZ = SLZ + PERC \qquad \text{(S29)}$$

$$Q_2 = SLZ \cdot k_2 \qquad \text{(S30)}$$

$$SLZ = SLZ - Q_2 \qquad \text{(S31)}$$

$$Q_t = Q_0 + Q_1 + Q_2 \qquad \text{(S32)}$$

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

$$Q_t^* = \int_0^{tmax} \xi(s:\theta_\alpha,\theta_\tau) \cdot Q(t-s)\,ds \qquad \text{(S33)}$$

$$\xi(s:\theta_\alpha,\theta_\tau) = \frac{1}{\Gamma(\theta_\alpha)\theta_\tau^{\theta_\alpha}} t^{\theta_\alpha-1} e^{-\frac{t}{\theta_\tau}} \qquad \text{(S34)}$$

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

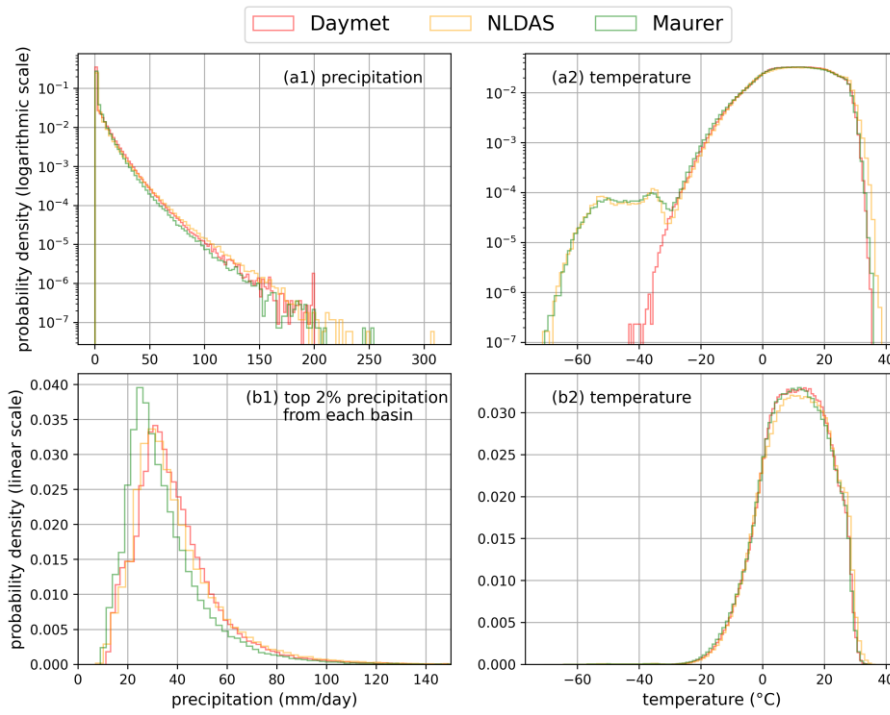*Figure B1. Probability density distributions (top panel in logarithmic scale, bottom panel in linear scale) of precipitation and temperature across three meteorological forcing datasets.*
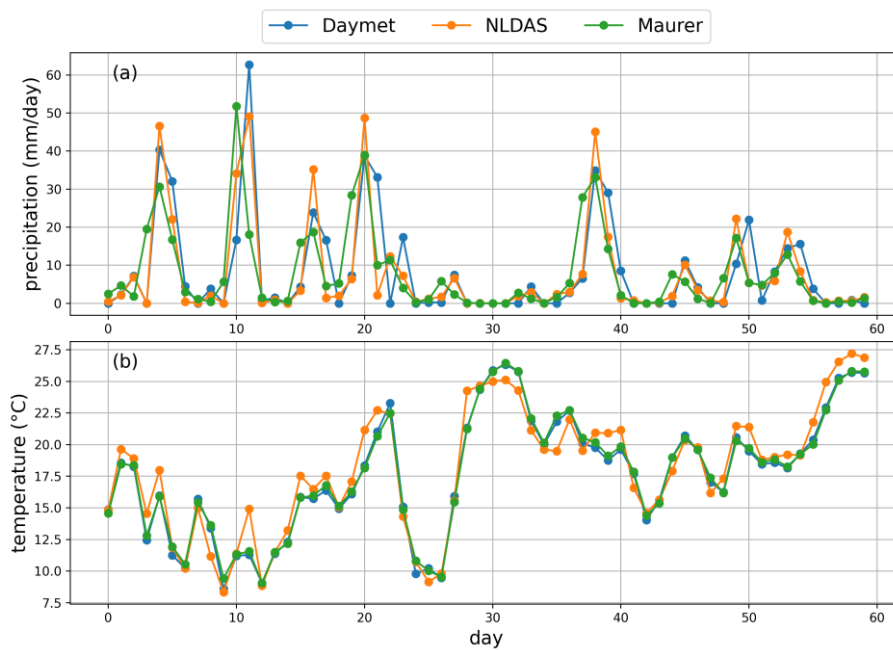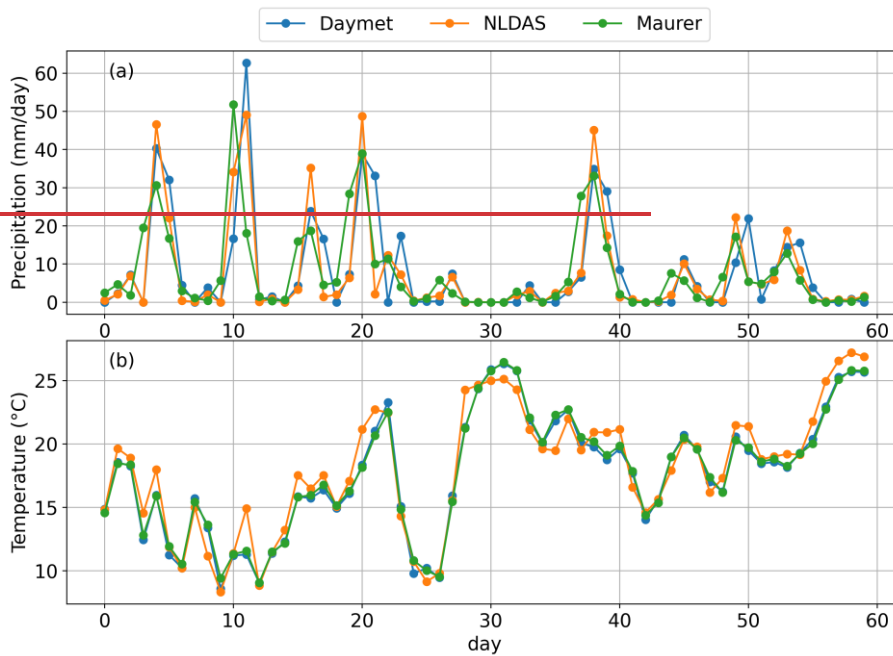
706



707

*Figure B2. ~~Illustrated~~Example of temporal variations ~~of~~in precipitation and temperature ~~in~~*

709     ~~a~~_for one_ basin across three meteorological forcing datasets.
710

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

712 *Table C1. Full names for the abbreviations of dynamic data (all but streamflow are "forcings")*

713 *and static basin attributes used as model inputs and outputs. All variables and their values are*

714 *provided in the CAMELS dataset (Addor et al., 2017) except for the NLDAS and Maurer daily*

715 *temperature extrema, which are from Kratzert et al. (2021). Potential evapotranspiration and*

716 *normalized streamflow were calculated in this work, using CAMELS data. The number in*

717 *parentheses ~~specifies~~indicates model usage: 1 ~~denotes use in~~for the LSTM model, and 2*

718 *~~denotes use in~~for the δHBV model.*

| Type | Abbreviation | Full name | Unit |
|---|---|---|---|
| **Dynamic data** | prcp (1,2) | Precipitation | mm/day |
| | pet (2) | Potential evapotranspiration (calculated in this work using the Hargreaves equation and CAMELS data) | mm/day |
| | tmean (2) | Mean air temperature | °C |
| | tmax (1) | Maximum air temperature | °C |
| | tmin (1) | Minimum air temperature | °C |
| | srad (1) | Shortwave radiation | $W/m^2$ |
| | vp (1) | Water vapor pressure | pa |
| | q_vol | Volumetric streamflow | ft3/s |
| | q (1,2) | Streamflow normalized by basin area (q_vol / area_gages2) | mm/day |
| **Static basin attributes** | p_mean (1,2) | Mean daily precipitation | mm/day |
| | pet_mean (1,2) | Mean daily potential evapotranspiration | mm/day |
| | p_seasonality (2) | Seasonality and timing of precipitation | - |
| | frac_snow (1,2) | Fraction of precipitation falling as snow | - |
| | aridity (1,2) | Rate of mean values of potential evapotranspiration and precipitation | - |
| | high_prec_freq (1,2) | Frequency of high precipitation days | days/year |
| | high_prec_dur (1,2) | Average duration of high precipitation events | days |
| | low_prec_freq (1,2) | 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 | $m^2$ |

719
720

721

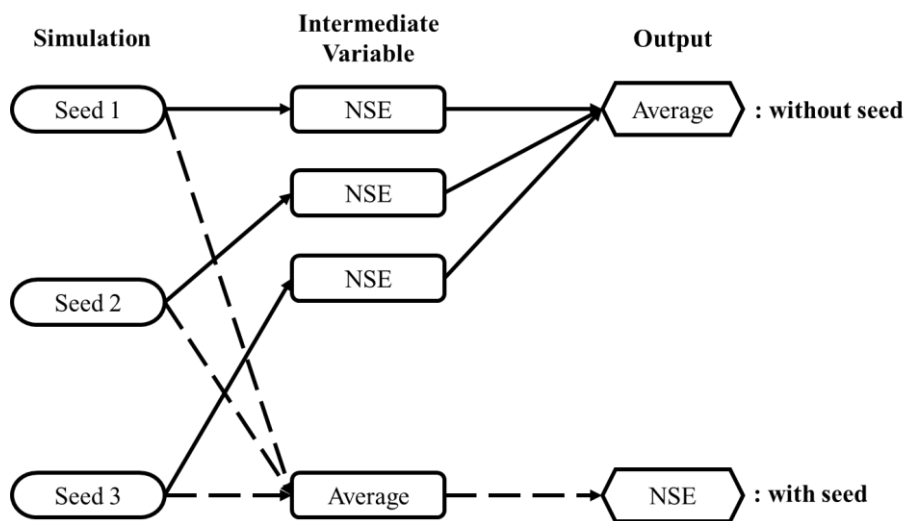*Figure C1. Ensemble frameworks to generate metrics for ensembles named without (solid arrows) and with (dashed arrows) "seed" as a subscript.*

724 *Table C2. ~~Evaluation~~Loss function and evaluation metrics.*

| Statistic | Equation* | Range | Optimal Value |
|---|---|---|---|
| Loss | $\dfrac{1}{n}\sum_{i=1}^{n}\dfrac{(O_i - S_i)^2}{(\sigma_O + \epsilon)^2}$ | 0.0 to $\infty$ | 0.0 |
| NSE | $1 - \dfrac{\sum_{i=1}^{n}(O_i - S_i)^2}{\sum_{i=1}^{n}(O_i - \mu_o)^2}$ | $-\infty$ to 1.0 | 1.0 |
| KGE | $1 - \sqrt{(r-1)^2 + (\beta-1)^2 + (\gamma-1)^2},$ $\beta = \dfrac{\mu_S}{\mu_O}, \gamma = \dfrac{cv_S}{cv_O} = \dfrac{\sigma_S/\mu_S}{\sigma_O/\mu_O}$ | $-\infty$ to 1.0 | 1.0 |
| PBIAS | $\dfrac{\sum_{i=1}^{n}(O_i - S_i)}{\sum_{i=1}^{n}O_i} \times 100$ | $-\infty$ to $\infty$ | 0.0 |
| RMSE | $\sqrt{\dfrac{1}{n}\sum_{i=1}^{n}(O_i - S_i)^2}$ | 0.0 to $\infty$ | 0.0 |
| spread | $\sqrt{\dfrac{1}{n}\dfrac{1}{e}\sum_{i=1}^{n}\sum_{j=1}^{e}(S_{i,j} - \mu_{S,i})^2}$ | 0.0 to $\infty$ | None |

725 * S is ~~a~~the streamflow simulation; O is the corresponding observation; n is the number of total

726 S or O; $\epsilon$ is a numerical stabilizer, with a default value of 0.1; e is the number of ensemble

727 members; r is the linear Pearson correlation between S and O; $\beta$ is the mean bias; and $\gamma$ is the

728 variability bias. The mean and standard deviation of simulations are denoted as $\mu_S$ and $\sigma_S$,

729 respectively, ~~and~~while $\mu_O$ and $\sigma_O$ ~~are the mean and standard deviation~~denote those of the

730 observations.

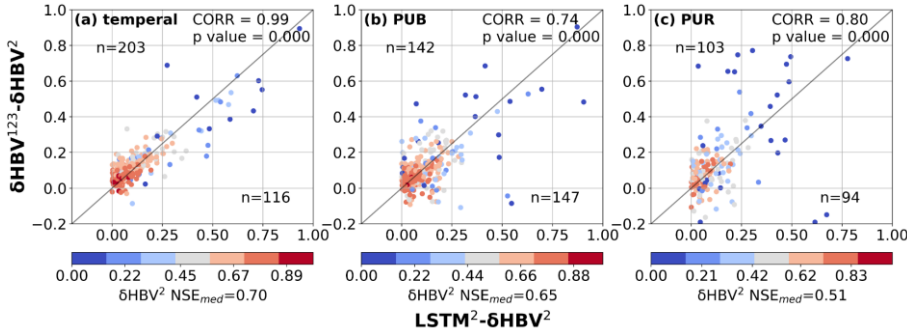**Appendix D: Additional details on model performance**

*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 NSE differences between $LSTM^2$ and $\delta HBV^2$ ($LSTM^2$ - $\delta HBV^2$), while the y-axis shows the NSE differences between $\delta HBV^{123}$ and $\delta HBV^2$ ($\delta HBV^{123}$- $\delta HBV^2$). Points are color-coded according to the NSE values of $\delta HBV^2$. The correlation coefficient (CORR) and p values between the x-axis values and the y-axis values, along with the median NSE value of $\delta HBV^2$ ($NSE_{med}$) on these basins, are also noted.*
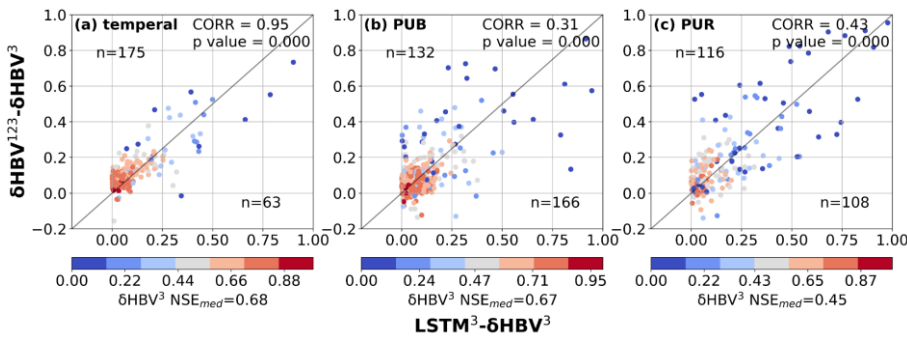


*Figure D2. Scatter plots comparing the performance differences between hydrological models for the basins where LSTM outperformed δHBV (the basins where δHBV outperformed are not shown in this plot). The x-axis represents the NSE differences between $LSTM^3$ and $\delta HBV^3$ ($LSTM^3$ - $\delta HBV^3$), while the y-axis shows the NSE differences between $\delta HBV^{123}$ and $\delta HBV^3$ ($\delta HBV^{123}$- $\delta HBV^3$). Points are color-coded according to the NSE values of $\delta HBV^3$. The correlation coefficient (CORR) and p values between the x-axis values and the y-axis values, along with the median NSE value of $\delta HBV^3$ ($NSE_{med}$) on these basins, are also noted.*
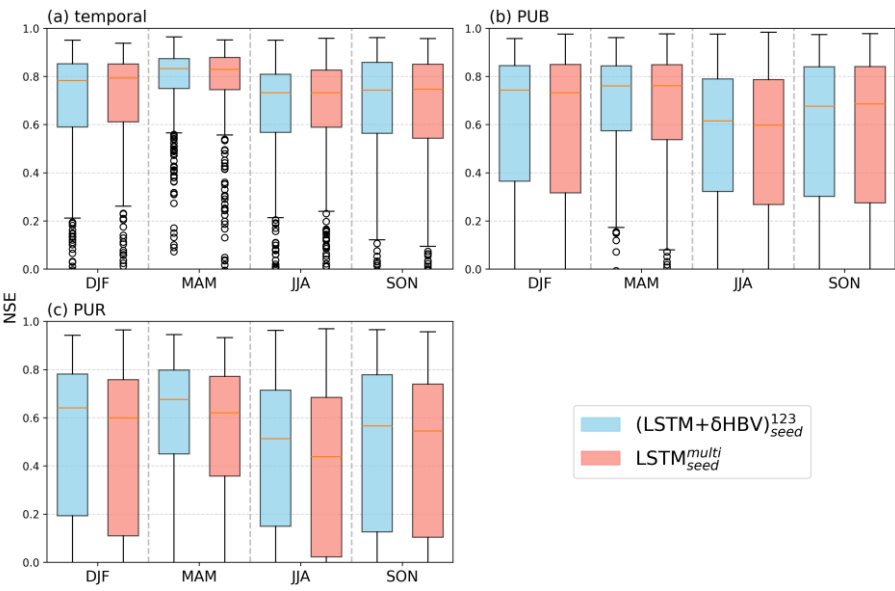
750



751

*Figure D3. Seasonal comparison of NSE values for $(LSTM + \delta HBV)^{123}_{seed}$ (blue) and $LSTM^{multi}_{seed}$ (red) in (a) temporal, (b) PUB, and (c) PUR tests. Each box represents the distribution of NSE values across 531 basins for a given season (DJF: December–February, MAM: March–May, JJA: June–August, SON: September–November). Vertical dashed lines separate different seasons. $(LSTM + \delta HBV)^{123}_{seed}$ performs better than $LSTM^{multi}_{seed}$ in most cases, especially during MAM, likely due to differences in snowmelt representation.*
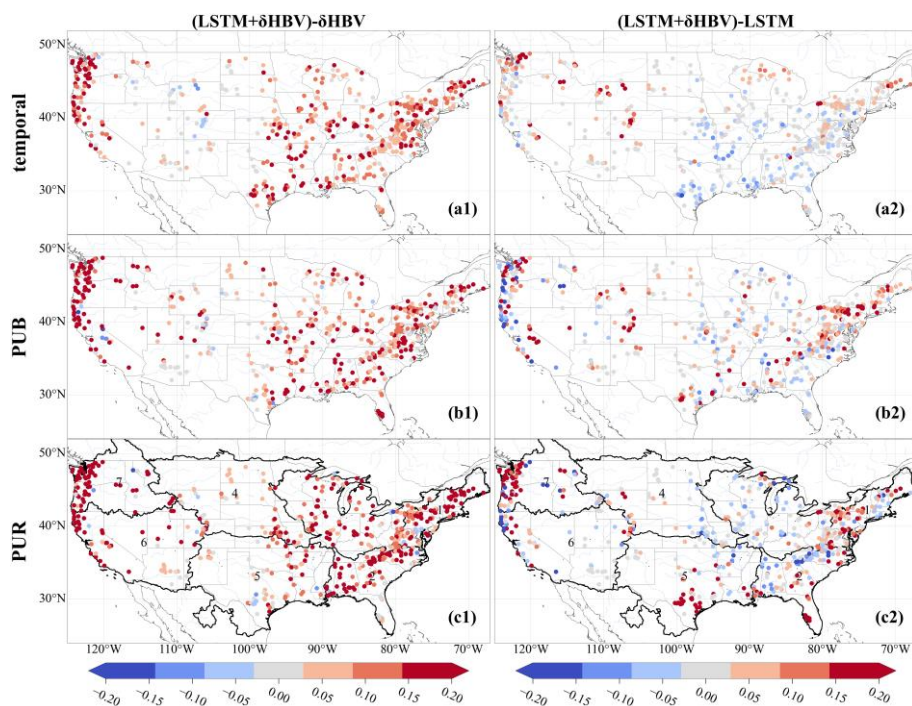
758

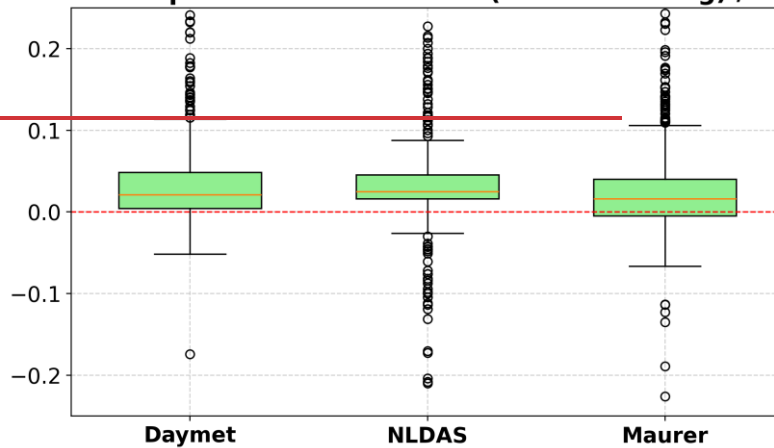759

760
761 *Figure ~~D3~~D4. Spatial distributions of model spread values increase from δHBV and LSTM to*
762 *the LSTM+δHBV ensemble across temporal, PUB, and PUR tests.*
763
764
765
766

Figure ~~D4~~_D5_. *Boxplot of relative temperature differences between the test and training periods,*
*calculated as* ~~(Test − Training) / Training~~_(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*
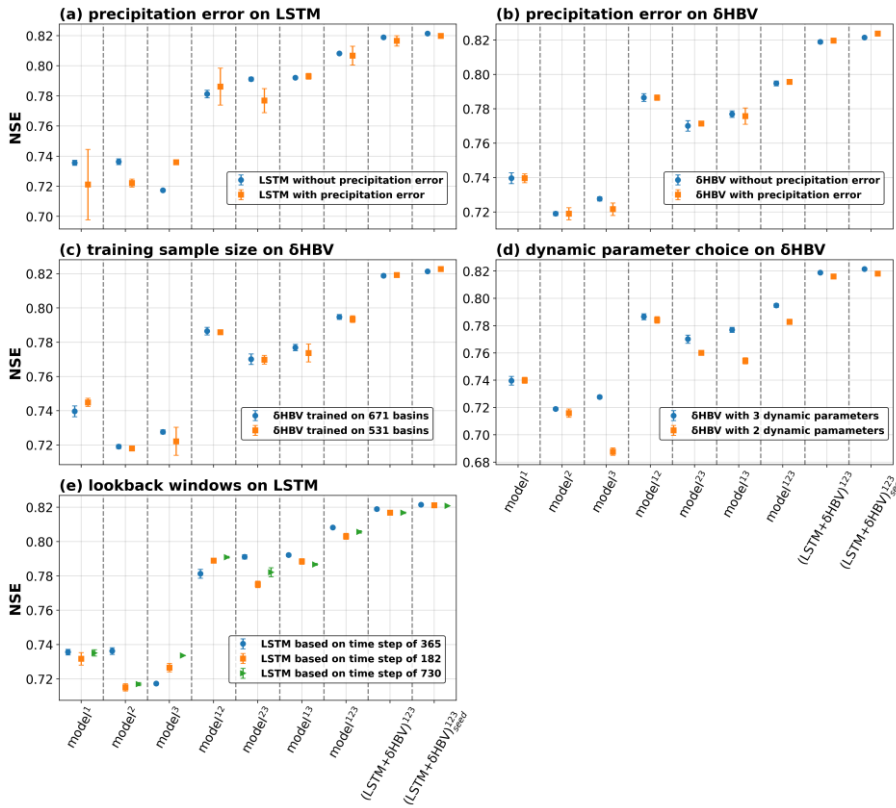*period relative to the training period.*

*Figure ~~D5~~D6. Simulation performance (NSE) under the temporal test: (a) LSTM model with and without a 10% precipitation error (precipitation × 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.*

789 *Table D1. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high*
790 *(highRMSE), and middle (midRMSE) flows based on 531 basins under the temporal test. The*
791 *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 |
| δHBV | NSE | 0.739688 | 0.71903 | 0.727669 |
| | KGE | 0.77033 | 0.730753 | 0.762022 |
| | RMSE | 1.18752 | 1.26239 | 1.23193 |
| | PBIAS | 5.07898 | -0.14449 | 3.65263 |
| | lowRMSE | 0.060906 | 0.063581 | 0.063466 |
| | highRMSE | 2.68479 | 3.13011 | 2.6845 |
| | midRMSE | 0.226595 | 0.245242 | 0.230125 |
| LSTM+δHBV | NSE | 0.787545 | 0.794053 | 0.790903 |

| | | | |
|---|---|---|---|
| KGE | 0.794412 | 0.78383 | 0.786854 |
| RMSE | 1.0777 | 1.0716 | 1.07141 |
| PBIAS | 4.59065 | 3.33053 | 3.45501 |
| lowRMSE | 0.059955 | 0.059565 | 0.054838 |
| highRMSE | 2.70216 | 2.88511 | 2.69633 |
| midRMSE | 0.20394 | 0.214726 | 0.212514 |

792
793
794

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

| Temporal | Number | Daymet+NLDAS | Daymet+Maurer | NLDAS+Maurer | All |
|---|---|---|---|---|---|
| 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 |
| δHBV | NSE | 0.786562 | 0.77012 | 0.776938 | 0.794796 |
| | KGE | 0.773732 | 0.778557 | 0.768854 | 0.77834 |
| | RMSE | 1.08362 | 1.12584 | 1.10875 | 1.06118 |
| | PBIAS | 1.91507 | 4.28194 | 2.03584 | 2.71021 |
| | lowRMSE | 0.061667 | 0.060679 | 0.062765 | 0.061539 |
| | highRMSE | 2.93961 | 2.7394 | 2.88758 | 2.84994 |
| | midRMSE | 0.230576 | 0.220743 | 0.230272 | 0.228375 |
| LSTM+δHBV | NSE | 0.811825 | 0.809964 | 0.811316 | 0.818907 |

| | | | | |
|---|---|---|---|---|
| KGE | 0.797564 | 0.797635 | 0.78735 | 0.794936 |
| RMSE | 1.01938 | 1.01755 | 1.0314 | 1.00067 |
| PBIAS | 4.14594 | 4.23333 | 3.19652 | 3.88096 |
| lowRMSE | 0.0603 | 0.058022 | 0.057882 | 0.059221 |
| highRMSE | 2.75275 | 2.67122 | 2.81393 | 2.70606 |
| midRMSE | 0.207637 | 0.205965 | 0.213191 | 0.207905 |

798

799  *Table D2. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high*
800  *(highRMSE), and middle (midRMSE) flows based on 531 basins under the PUB test. The values*
801  *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 |

| | | | |
|---|---|---|---|
| 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 |

802

803

804

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

| PUB | Number | Daymet+NLDAS | Daymet+Maurer | NLDAS+Maurer | All |
|---|---|---|---|---|---|
| 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 |
| δHBV | 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+δHBV | NSE | 0.780625 | 0.764866 | 0.767761 | 0.785833 |

| | | | | |
|---|---|---|---|---|
| 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 |

808
809

810 *Table D3. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high*
811 *(highRMSE), and middle (midRMSE) flows based on 531 basins under the PUR test. The values*
812 *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 |

| | | | |
|---|---|---|---|
| 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 |

813
814

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

| PUR | Number | Daymet+NLDAS | Daymet+Maurer | NLDAS+Maurer | All |
|---|---|---|---|---|---|
| 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 |
| δHBV | 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+δHBV | NSE | 0.685032 | 0.680872 | 0.679321 | 0.700814 |

| | | | | |
|---|---|---|---|---|
| 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 |

818
819

820  *Table D4. Median NSE, KGE, RMSE, PBIAS, and RMSE values under low (lowRMSE), high*
821  *(highRMSE), and middle (midRMSE) flows based on 531 basins under the temporal, PUB, and*
822  *PUR tests of $LSTM^{multi}$, $(LSTM + \delta HBV)^{123} + LSTM^{multi}$, their "seed" version, and*
823  *$(LSTM + \delta HBV)^{123}_{seed}$.*

| 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 |

| 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 |

824

825   *Table D4 (continued). Median NSE, KGE, RMSE, PBIAS, and RMSE values under low*
826   *(lowRMSE), high (highRMSE), and middle (midRMSE) flows based on 531 basins under the*
827   *temporal, PUB, and PUR tests of $LSTM^{multi}$, $(LSTM + \delta HBV)^{123} + LSTM^{multi}$ , their "seed"*
828   *version, and $(LSTM + \delta HBV)^{123}_{seed}$.*

| 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 |

| | | | | |
|---|---|---|---|---|
| | NSE | 0.705154 | 0.665723 | 0.701504 |
| | KGE | 0.651538 | 0.614649 | 0.64373 |
| | RMSE | 1.30377 | 1.3727 | 1.2851 |
| PUR | 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 |

829

830 *Table D5. Median NSE values based on ten different random seeds during the temporal test.*
831 *Each number (1 through 10) represents metric values calculated for an individual simulation*
832 *based on only one random seed. "Seed" indicates metric values calculated by averages of these*
833 *ten simulations based on different random seeds, while "mean" denotes the average of metrics*
834 *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 |

835

**Appendix E: Intuitive visualization of the relative contributions of ensemble members**

**based on optimized weights**



*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 δHBV 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.*

853
854 *Figure E3. Spatial distributions of weights of the Daymet, NLDAS, and Maurer meteorological*
855 *forcing datasets, estimated by a genetic algorithm based on streamflow observations during*
856 *the test periods. The weights are normalized by the maximum weight within each ensemble*
857 *group. These weights are used exclusively for qualitatively analyzing the relative contributions*
858 *of different ensemble members, with higher values indicating larger relative contributions.*
859

860

861

862 *Table E1. Comparisons of metric values between averaged ensemble simulations and*
863 *optimized weighted simulations, estimated using a genetic algorithm based on streamflow*
864 *observations during the test periods. The results highlight the potential for further*
865 *improvements in ensemble simulations.*
866

|  | Temporal | Averaged | Optimized weighted |
|---|---|---|---|
| Temporal | NSE | 0.821444 | 0.844303212 |
|  | KGE | 0.795317 | 0.829996445 |
|  | RMSE | 0.99455 | 0.920954559 |
|  | PBIAS | 3.99009 | 3.252278013 |
|  | lowRMSE | 0.059782 | 0.057137161 |
|  | highRMSE | 2.7279 | 2.451194907 |
|  | midRMSE | 0.209943 | 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 |
|  | 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 |

867

868

**Code and data availability**

The source codes and datasets utilized in this study are publicly accessible through the following repositories: The δHBV modeling framework, including all computational scripts and documentation, is hosted on Zenodo (https://doi.org/10.5281/zenodo.7091334) ~~(Feng et al., 2023a)~~(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)~~(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)~~(Addor et al., 2017; Newman and Clark, 2014). The streamflow simulations produced in this study will be made available on Zenodo upon acceptance of the manuscript.

**Author contributions**

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

901

**References**

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

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

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

912 Bandai, T. and Ghezzehei, T. A.: Physics-informed neural networks with monotonicity
913 constraints for Richardson-Richards equation: Estimation of constitutive relationships and soil
914 water flux density from volumetric water content measurements, Water Resources Research,
915 57, e2020WR027642, https://doi.org/10.1029/2020wr027642, 2021.

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

919 Beck, H. E., Pan, M., Lin, P., Seibert, J., Dijk, A. I. J. M. van, and Wood, E. F.: Global fully
920 distributed parameter regionalization based on observed streamflow from 4,229 headwater
921 catchments, Journal of Geophysical Research: Atmospheres, 125, e2019JD031485,
922 https://doi.org/10.1029/2019JD031485, 2020.

923 Behnke, R., Vavrus, S., Allstadt, A., Albright, T., Thogmartin, W. E., and Radeloff, V. C.:
924 Evaluation of downscaled, gridded climate data for the conterminous United States, Ecological
925 Applications, 26, 1338–1351, https://doi.org/10.1002/15-1061, 2016.

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

Bellmore, J. R., 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, Wiley Interdisciplinary Reviews: Water, 4, e1164, https://doi.org/10.1002/wat2.1164, 2017.

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-

969 based models with multiphysical outputs can approach state-of-the-art hydrologic prediction
970 accuracy, Water Resources Research, 58, e2022WR032404,
971 https://doi.org/10.1029/2022WR032404, 2022.

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

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

978 Frame, J. M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., Qualls, L. M., Gupta,
979 H. V., and Nearing, G. S.: Deep learning rainfall–runoff predictions of extreme events,
980 Hydrology and Earth System Sciences, 26, 3377–3392, https://doi.org/10.5194/hess-26-3377-
981 2022, 2022.

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

985 Hargreaves, G. H.: Defining and using reference evapotranspiration, Journal of Irrigation and
986 Drainage Engineering, 120, 1132–1139, https://doi.org/10.1061/(ASCE)0733-
987 9437(1994)120:6(1132), 1994.

988 He, Y., Chen, M., Wen, Y., Duan, Q., Yue, S., Zhang, J., Li, W., Sun, R., Zhang, Z., Tao, R.,
989 Tang, W., and Lü, G.: An open online simulation strategy for hydrological ensemble
990 forecasting, Environmental Modelling & Software, 174, 105975,
991 https://doi.org/10.1016/j.envsoft.2024.105975, 2024.

992 Heidari, H., Arabi, M., Warziniack, T., and Kao, S.-C.: Assessing shifts in regional
993 hydroclimatic conditions of U.S. river basins in response to climate change over the 21st
994 century, Earth's Future, 8, e2020EF001657, https://doi.org/10.1029/2020EF001657, 2020.

995 Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, Neural Computation, 9, 1735–
996 1780, https://doi.org/10.1162/neco.1997.9.8.1735, 1997.

997 Ji, H., Song, Y., Bindas, T., Shen, C., Yang, Y., Pan, M., Liu, J., Rahmani, F., Abbas, A., Beck,
998 H., Lawson, K., and Wada, Y.: Distinct hydrologic response patterns and trends worldwide
999 revealed by physics-embedded learning, Nat. Commun., 16, 9169,
1000 https://doi.org/10.1038/s41467-025-64367-1, 2025.

1001 Jiang, S., Zheng, Y., and Solomatine, D.: Improving AI system awareness of geoscience
1002 knowledge: Symbiotic integration of physical approaches and deep learning, Geophys. Res.
1003 Lett., 47, e2020GL088229, https://doi.org/10.1029/2020GL088229, 2020.

1004 Kling, H., Fuchs, M., and Paulin, M.: Runoff conditions in the upper Danube basin under an
1005 ensemble of climate change scenarios, Journal of Hydrology, 424–425, 264–277,
1006 https://doi.org/10.1016/j.jhydrol.2012.01.011, 2012.

1007 Kraft, B., Jung, M., Körner, M., Koirala, S., and Reichstein, M.: Towards hybrid modeling of
1008 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/gg4ck8, 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., Shi, L., and Zhong, H.: Assessing the Global Relationships Between Teleconnection Factors and Terrestrial Water Storage Components, Water Resources Management, 36, 119–133, https://doi.org/10.1007/s11269-021-03015-x, 2022.

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.

77

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.

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, 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, 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.

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, 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, 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.

1171 Yilmaz, K. K., Gupta, H. V., and Wagener, T.: A process-based diagnostic approach to model
1172 evaluation: Application to the NWS distributed hydrologic model, Water Resources Research,
1173 44, https://doi.org/10/fpvsgb, 2008.

1174 Yu, D., Yang, J., Shi, L., Zhang, Q., Huang, K., Fang, Y., and Zha, Y.: On the uncertainty of
1175 initial condition and initialization approaches in variably saturated flow modeling, Hydrology
1176 and Earth System Sciences, 23, 2897–2914, https://doi.org/10.5194/hess-23-2897-2019, 2019.

1177 Yu, M., Huang, Q., and Li, Z.: Deep learning for spatiotemporal forecasting in Earth system
1178 science: a review, International Journal of Digital Earth, 17, 2391952,
1179 https://doi.org/10.1080/17538947.2024.2391952, 2024.

1180 Zhang, Q., Shi, L., Holzman, M., Ye, M., Wang, Y., Carmona, F., and Zha, Y.: A dynamic
1181 data-driven method for dealing with model structural error in soil moisture data assimilation,
1182 Advances in Water Resources, 132, 103407, https://doi.org/10.1016/j.advwatres.2019.103407,
1183 2019.

1184 Zounemat-Kermani, M., Batelaan, O., Fadaee, M., and Hinkelmann, R.: Ensemble machine
1185 learning paradigms in hydrology: A review, Journal of Hydrology, 598, 126266,
1186 https://doi.org/10.1016/j.jhydrol.2021.126266, 2021.

1187

1188