1 **On the need for physical constraints in deep leaning rainfall-runoff**

2 **projections under climate change: a sensitivity analysis to warming and shifts**

3 **in potential evapotranspiration**

4

5 **Sungwook Wi[1], Scott Steinschneider[1]**

6 [1]Department of Biological and Environmental Engineering, Cornell University, Ithaca, NY, USA

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

**Abstract**

Deep learning (DL) rainfall-runoff models outperform process-based models in a range of applications. However, it remains unclear whether DL models can produce physically plausible projections of streamflow under climate change. We investigate this question through a sensitivity analysis of modeled responses to increases in temperature and potential evapotranspiration (PET), with other meteorological variables left unchanged. Previous research has shown that temperature-based PET methods overestimate evaporative water loss under warming compared to energy budget-based PET methods. We therefore assume that reliable streamflow responses to warming should exhibit less evaporative water loss when forced with smaller, energy budget-based PET compared to temperature-based PET. We conduct this assessment using three conceptual, process-based rainfall-runoff models and three DL models, trained and tested across 212 watersheds in the Great Lakes basin. The DL models include a Long Short-Term Memory network (LSTM), a mass-conserving LSTM (MC-LSTM), and a novel variant of the MC-LSTM that also respects the relationship between PET and evaporative water loss (MC-LSTM-PET). After validating models against historical streamflow and actual evapotranspiration, we force all models with scenarios of warming, historical precipitation, and both temperature-based (Hamon) and energy budget-based (Priestley-Taylor) PET, and compare their responses in long-term mean daily flow, low flows, high flows, and seasonal streamflow timing. We also explore similar responses using a National LSTM fit to 531 watersheds across the United States to assess how the inclusion of a larger and more diverse set of basins influences signals of hydrologic response under warming. The main results of this study are as follows:

1. The three Great Lakes DL models substantially outperform all process models in streamflow estimation. The MC-LSTM-PET also matches the best process models and outperforms the MC-LSTM in estimating actual evapotranspiration.

2. All process models show a downward shift in long-term mean daily flows under warming, but median shifts are considerably larger under temperature-based PET (-17% to -25%) than energy budget-based PET (-6% to -9%). The MC-LSTM-PET model exhibits similar differences in water

50         loss across the different PET forcings. Conversely, the LSTM exhibits unrealistically large water

51         losses under warming using Priestley-Taylor PET (-20%), while the MC-LSTM is relatively

52         insensitive to PET method.

53    3.  DL models exhibit smaller changes in high flows and seasonal timing of flows as compared to the

54        process models, while DL estimates of low flows are within the range estimated by the process

55        models.

56    4.  Like the Great Lakes LSTM, the National LSTM also shows unrealistically large water losses under

57        warming (-25%), but it is more stable when many inputs are changed under warming and better

58        aligns with process model responses for seasonal timing of flows.

59    Ultimately, the results of this sensitivity analysis suggest that physical considerations regarding model

60    architecture and input variables may be necessary to promote the physical realism of deep learning-based

61    hydrologic projections under climate change.

62

63    **Keywords**

64    Deep learning, machine learning, Long Short-Term Memory network, LSTM, Great Lakes, climate

65    change, rainfall-runoff

66

67

68

69

70

71

72

73

74

## 1. Introduction

Rainfall-runoff models are used throughout hydrology in a range of applications, including retrospective streamflow estimation (Hansen et al. 2019), streamflow forecasting (Demargne et al., 2014), and prediction in ungauged basins (Hrachowitz et al., 2013). Work over the last few years has demonstrated that deep learning (DL) rainfall-runoff models (e.g., Long Short-Term Memory networks (LSTMs); Hochreiter and Schmidhuber, 1997) outperform conventional process-based models in each of these applications, especially when those DL models are trained with large datasets collected across watersheds with diverse climates and landscapes (Kratzert et al., 2019a,b; Feng et al., 2020; Ma et al., 2021; Gauch et al., 2021a,b; Nearing et al., 2021). For example, in one extensive benchmarking study, Mai et al. (2022) found that a regionally trained LSTM outperformed 12 other lumped and distributed process-based models of varying complexity in rivers and streams throughout the Great Lakes basin. These and similar results have led some to argue that DL models represent the most accurate and extrapolatable rainfall-runoff models available (Nearing et al., 2022).

However, there remains one use case of rainfall-runoff models where the superiority of DL is unclear: long-term projections of streamflow under climate change. Past studies using DL rainfall-runoff models for hydrologic projections under climate change are rare (Lee et al., 2020; Li et al., 2022), and few have evaluated their physical plausibility (Razavi, 2021; Reichert et al., 2023; Zhong et al., 2023). A reasonable concern is whether DL rainfall-runoff models can extrapolate hydrologic response under unprecedented climate conditions, given that they are entirely data driven and do not explicitly represent the physics of the system. It is not clear *a priori* whether this concern has merit, because DL models fit to a large and diverse set of basins have the benefit of learning hydrologic response across climate and landscape gradients. In so doing, the model can, for example, learn hydrologic responses to climate in warmer regions and then transfer this knowledge to projections of streamflow in cooler regions subject to climate change induced warming. In addition, past work has shown that LSTMs trained only to predict streamflow have memory cells that strongly correlate with independent measures of soil moisture and snowpack (Lees et al. 2022),

4

101    suggesting that DL hydrologic models can learn fundamental hydrologic processes. A potential implication

102    of this finding might be that these models can produce physically plausible streamflow predictions under

103    new climate conditions.

104

105    It is challenging to assess the physical plausibility of DL-based hydrologic projections under substantially

106    different climate conditions, because there are no future observations against which to compare. This

107    challenge is exacerbated by significant uncertainty in process model projections under alternative climates,

108    which makes establishing reliable benchmarks difficult. Future process model-based projections can vary

109    widely due to both parametric and structural uncertainty (Bastola et al., 2011; Clark et al., 2016; Melsen et

110    al., 2018), and even for models that exhibit similar performance under historical conditions (Krysanova et

111    al., 2018). Assumptions around stationary model parameters are not always valid (Merz et al., 2011;

112    Wallner and Haberlandt, 2015), and added complexity for improved process representation is not always

113    well supported by data (Clark et al., 2017; Towler et al., 2023; Yan et al., 2023). Together, these challenges

114    highlight the difficulty in establishing good benchmarks of hydrologic response under alternative climates

115    against which to compare and evaluate DL-based hydrologic projections under climate change.

116

117    Recently, Wi and Steinschneider (2022) (hereafter WS22) forwarded an experimental design to evaluate

118    the physical plausibility of DL hydrologic responses to new climates, in which DL hydrologic models were

119    forced with historical precipitation and temperature, but with temperatures adjusted by up to 4°C. Based on

120    past literature, WS22 posited that in non-glaciated regions, physically plausible hydrologic responses

121    should show an increase in water loss, defined as water that enters the watershed via precipitation but never

122    contributes to streamflow because it is 'lost' to a terminal sink. Specifically, WS22 assumed that

123    evaporative water loss should increase and annual average streamflow should decline compared to a

124    baseline simulation due to increases in potential evapotranspiration (PET) with warming (and no changes

125    in precipitation). Results showed that an LSTM trained to the 15 watersheds in California often led to

126    misleading increases in annual runoff under warming, while this phenomenon was less likely (though still

5

127 present) in a DL model trained to 531 catchments across the United States. WS22 also conducted their

128 experiment with physics-informed machine learning (PIML) models (Karpatne et al., 2017), using process

129 model output directly as input to the LSTM (similar to Konapala et al., 2020; Lu et al., 2021; Frame et al.,

130 2021a) or as additional target variables in a multi-output architecture. The former approach had some

131 success in removing instances of increasing runoff ratio with warming, although this was dependent on the

132 process model used.

133

134 Other PIML approaches that more directly adjust the architecture of DL rainfall-runoff models may be

135 better suited for improving long-term streamflow projections under climate change without requiring an

136 accurate process-based model. For instance, Hoedt et al. (2021) introduced a mass conserving LSTM (MC-

137 LSTM) that ensures cumulative streamflow predictions do not exceed precipitation inputs. Hybrid models

138 present a related approach, where DL modules are embedded within process models structures (Jiang et al.,

139 2020; Feng et al., 2022; Hoge et al., 2022; Feng et al., 2023a). In some cases, these architectural changes

140 can degrade performance compared to a standard LSTM (Frame et al., 2021b; Feng et al., 2023b), but other

141 times such changes can be beneficial (Feng et al., 2023a). Some have argued that these physical constraints

142 may inhibit the ability of DL models to learn biases in forcing data (Frame et al. 2022), but the benefits of

143 such mass conserving architectures have not been tested when employed under previously unobserved

144 climate change.

145

146 For all models considered in WS22, a major focus was evaluating the direction of annual total runoff change

147 in the presence of warming and no change in precipitation. However, that study did not consider the

148 magnitude of runoff change and how it relates to projected changes in PET. As we argue below, this

149 comparison provides a unique way to assess the physical plausibility of future hydrologic projections.

150 Several studies have investigated the effects of different PET estimation methods on the magnitude of PET

151 and runoff change in a warming climate (Lofgren et al., 2011; Shaw and Riha, 2011; Lofgren and Rouhana,

152 2016; Milly and Dunne, 2017; Lemaitre-Basset et al. 2022). Broadly, these studies have shown that

153  temperature-based PET estimation methods (e.g., Hamon, Thornthwaite) substantially overestimate

154  increases in PET under warming as compared to energy budget-based PET estimation methods (e.g.,

155  Penman-Monteith, Priestley-Taylor), and consequently lead to unrealistic declines in streamflow under

156  climate change. This is because the actual drying power of the atmosphere is driven by the availability of

157  energy at the surface from net radiation, the current moisture content of the air, temperature (and its effect

158  on the water holding capacity of the air and vapor pressure deficit), and wind speeds. Energy budget-based

159  methods, while imperfect and at times empirical (Greve et al. 2019; Liu et al., 2022), account for some or

160  all of these factors in ways that are generally consistent with their causal impact on PET, while temperature-

161  based methods estimate PET using strictly empirical relationships based largely or entirely on temperature.

162  The latter approach works sufficiently well for rainfall-runoff modeling under historical conditions because

163  of the strong correlation between temperature, net radiation, and PET on seasonal timescales, even though

164  this correlation weakens considerably at shorter timescales (Lofgren et al., 2011). Under climate change,

165  consistent and prominent increases are projected for temperature, but projected changes are less prominent

166  or more uncertain for other factors affecting PET (Lin et al., 2018; Pryor et al., 2020, Liu et al. 2020).

167  Consequently, temperature-based PET methods substantially overestimate future projections of PET

168  compared to energy budget-based methods (Lofgren et al., 2011; Shaw and Riha, 2011; Lofgren and

169  Rouhana, 2016; Milly and Dunne, 2017; Lemaitre-Basset et al. 2022).

170

171  As argued by Lofgren and Rouhana (2016), the bias in PET and runoff that results from different PET

172  estimation methods under warming provides a unique opportunity to assess the physically plausibility of

173  hydrologic projections under climate change. In this study, we adopt this strategy for DL rainfall-runoff

174  models through a sensitivity analysis in which both conceptual, process-based and DL hydrologic models

175  are trained with either temperature-based or energy budget-based estimates of PET, along with other

176  meteorological data (precipitation, temperature). These models are then forced with the historical

177  precipitation and temperature series, but with the temperatures warmed by an additive factor and PET

178  calculated from the warmed temperatures using both PET estimation methods. We show that the process

179   models 1) exhibit similar performance in historical training and testing periods when using either

180   temperature-based or energy budget-based PET estimates; but 2) exhibit substantially larger long-term

181   mean streamflow declines under warming when using future PET estimated with a temperature-based

182   method. If the DL rainfall-runoff models follow the same pattern, this would suggest that these models are

183   able to learn the role of PET on evaporative water loss. However, if DL-based models estimate similarly

184   large long-term mean streamflow declines regardless of the method used to estimate and project PET, this

185   would suggest that the DL models did not learn a mapping between PET and evaporative water loss. Rather,

186   the DL models learned the historical (but non-causal) correlation between temperature and evaporative

187   water loss, and then incorrectly extrapolated that effect into the future with warmer temperatures. We show

188   this latter outcome to be the case, which indicates that some degree of PIML may be necessary to guide a

189   DL model towards physically plausible projections under climate change.

190

191   We conduct the experiment above in a case study on 212 watersheds across the Great Lakes basin, using

192   both standard and PIML-based LSTMs. We show that a standard LSTM produces unrealistic hydrologic

193   responses to warming because it relies on historical and geographically pervasive correlations between

194   temperature and PET to estimate streamflow losses. We also show that PIML-based DL models are better

195   able to relate changes in temperature and PET to streamflow change, especially those PIML approaches

196   that directly map PET to evaporative water loss in their architecture.

197

198   The primary goal of this work is to forward an experimental design that can be used to evaluate the

199   suitability of DL rainfall-runoff models for hydrologic projections under climate change, in line with a

200   recent call to design benchmarking studies that assess whether models are fit for specific purposes (Beven,

201   2023). The Great Lakes provides an important case study for this work, given their importance to the culture,

202   ecosystems, and economy of North America (Campbell et al., 2015; Steinman et al., 2017). Projections of

203   future water supplies and water levels in the Great Lakes are highly uncertain (Gronewold and Rood, 2019),

204   in part because of uncertainty in future runoff draining into the lakes from a large contributing area
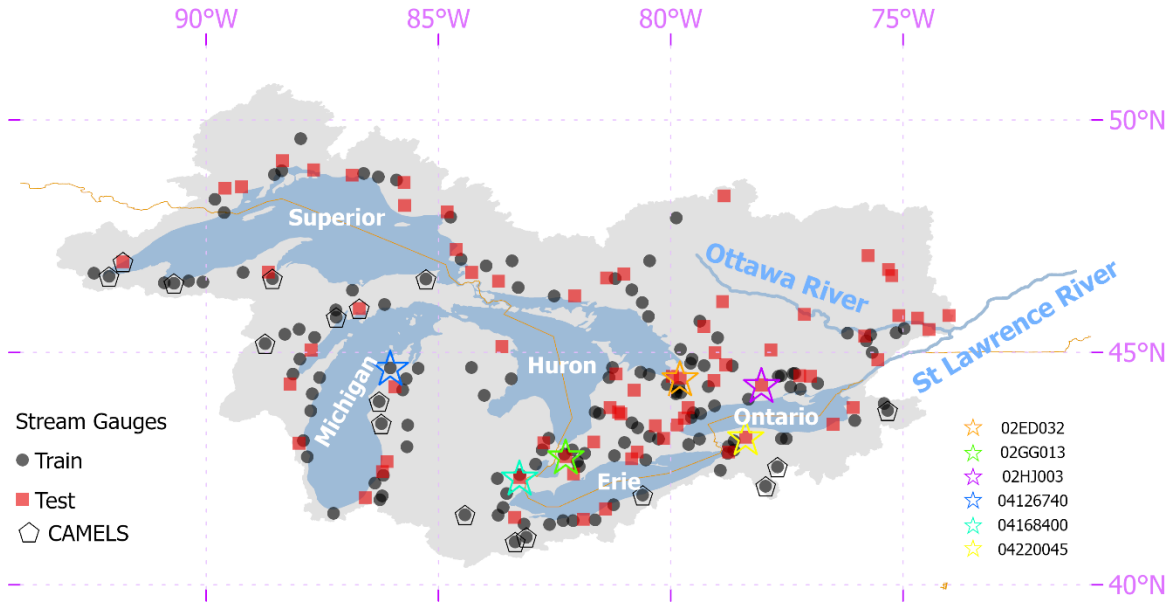
205 (Kayastha et al. 2022), much of which is ungauged (Fry et al., 2013). Improved rainfall-runoff models that

206 can regionalize across the entire Great Lakes basin are necessary to help address this challenge, and so an

207 auxiliary goal of this work is to contribute PIML rainfall-runoff models to the Great Lakes Runoff

208 Intercomparison Project Phase 4 presented in Mai et al. (2022). This study currently provides one of the

209 most robust benchmarks comparing DL rainfall-runoff models to a range of process-based models, and so

210 we design our experiment to be consistent with the data and model development rules outlined in that

211 intercomparison project.

212

213 **2. Data**

214 This study focuses on 212 watersheds draining into the Great Lakes and Ottawa River, which are all located

215 in the St. Lawrence River basin (Figure 1). For direct comparability to previous results from the Great Lakes

216 Runoff Intercomparison Project, all data for these watersheds are taken directly from the work in Mai et al.

217 (2022) and include daily streamflow time series, meteorological forcings, geophysical attributes for each

218 watershed, and auxiliary hydrologic fluxes. Daily streamflow were gathered from the U.S. Geological

219 Survey and Water Survey Canada between January 2000 and December 2017. All streamflow gauging

220 stations have a drainage area greater than or equal to 200 km$^2$ and less than 5% missing data in the study

221 period. The watersheds are evenly distributed across the five lake basins and the Ottawa River basin, and

222 they represent a range of land use/land cover types and degrees of hydrologic alteration from human activity.

223 In the experiments described further below, 141 of the watersheds are designated as training sites, and the

224 remaining 71 watersheds are used for testing (see Figure 1). In addition, the period between January 2000

225 to December 2010 is reserved for model training (termed the training period), and the period between

226 January 2011 – December 2017 is used for model testing (termed the testing period).

227

**Figure 1.** Great Lakes domain, with training and testing streamflow gauges used throughout this study. A subset of seventeen of these gauges that are also in the CAMELS database are highlighted, as are six sites used to present select results in Section 4.

Meteorological forcings are taken from the Regional Deterministic Reanalysis System v2, which is an hourly, 10 km dataset available across North America (Gasset et al., 2021). Hourly precipitation, net incoming shortwave radiation ($R_s$), and temperature are aggregated into a basin-wide daily precipitation average, daily $R_s$ average, and daily minimum and maximum temperature. We note that the precipitation data from the Regional Deterministic Reanalysis System v2 is produced from the Canadian Precipitation Analysis, which combines available surface observations of precipitation with a short-term reforecast provided by the 10 km Regional Deterministic Reforecast System. That is, the precipitation data is not model based, but rather is based on gauged data and spatially interpolated using information from modeled output.

Geophysical attributes for each watershed were collected from a variety of sources. Basin-average statistics of elevation and slope were derived from the HydroSHEDS dataset (Lehner et al., 2008), which provides a

10

245 digital elevation model with 3 arcsec resolution. Soil properties (e.g., soil texture, classes) were gathered

246 from the Global Soil Dataset for Earth System Models (Shangguan et al., 2014), which is available at a 30

247 arcsec resolution. Land cover data at a 30 m resolution and based on Landsat imagery from 2010-2011 were

248 derived from the North American Land Change Monitoring System (NALCMS, 2017). These geophysical

249 datasets were used to derive basin-averaged attributes for each watershed, listed in Table 1.

250

251 **Table 1**. Watershed attributes used in the deep learning models developed in this work (adapted from Mai
252 et al., 2022).

| Attribute | Description |
|---|---|
| p_mean | Mean daily precipitation |
| pet_mean | Mean daily potential evapotranspiration |
| aridity | Ratio of mean PET to mean precipitation |
| t_mean | Mean of daily maximum and daily minimum temperature |
| frac_snow | Fraction of precipitation falling on days with mean daily temperatures below 0°C |
| high_prec_freq | Fraction of high-precipitation days (= 5 times mean daily precipitation) |
| high_prec_dur | Average duration of high-precipitation events (number of consecutive days with = 5 times mean daily precipitation) |
| low_prec_freq | Fraction of dry days (< 1 mm d-1 daily precipitation) |
| low_prec_dur | Average duration of dry periods (number of consecutive days with daily precipitation < 1 mm d-1) |
| mean_elev | Catchment mean elevation |
| std_elev | Standard deviation of catchment elevation |
| mean_slope | Catchment mean slope |
| std_slope | Standard deviation of catchment slope |
| area_km2 | Catchment area |
| Temperate-or-sub-polar-needleleaf-forest | Fraction of land covered by "Temperate-or-sub-polar-needleleaf-forest" |
| Temperate-or-sub-polar-grassland | Fraction of land covered by "Temperate-or-sub-polar-grassland" |
| Temperate-or-sub-polar-shrubland | Fraction of land covered by "Temperate-or-sub-polar-shrubland" |
| Temperate-or-sub-polar-grassland | Fraction of land covered by "Temperate-or-sub-polar-grassland" |
| Mixed-Forest | Fraction of land covered by "Mixed-Forest" |
| Wetland | Fraction of land covered by "Wetland" |

| Cropland | Fraction of land covered by "Cropland" |
|---|---|
| Barren-Lands | Fraction of land covered by "Barren-Lands" |
| Urban-and-Built-up | Fraction of land covered by "Urban-and-Built-up" |
| Water | Fraction of land covered by "Water" |
| BD | Soil bulk density (g cm-3) |
| CLAY | Soil clay content (% of weight) |
| GRAV | Soil gravel content (% of volume) |
| OC | Soil organic carbon (% of weight) |
| SAND | Soil sand content (% of weight) |
| SILT | Soil silt content (% of weight) |

253

254 Finally, we also collect daily actual evapotranspiration (AET) for each watershed in millimeters per day,

255 which was originally taken from the Global Land Evaporation Amsterdam Model (GLEAM) v3.5b dataset

256 (Martens et al., 2017). GLEAM couples remotely sensed observations of microwave Vegetation Optical

257 Depth, a multi-layer soil moisture model driven by observed precipitation and assimilating satellite surface

258 soil moisture observations, and Priestly-Taylor based estimates of PET to derive an estimate of AET for

259 each day. The daily data were originally available over the entire study domain at a 0.25° resolution between

260 2003-2017 and were aggregated to basin-wide totals for each watershed. While AET from GLEAM is still

261 uncertain, it provides a useful, independent, remote-sensing based benchmark against which to compare
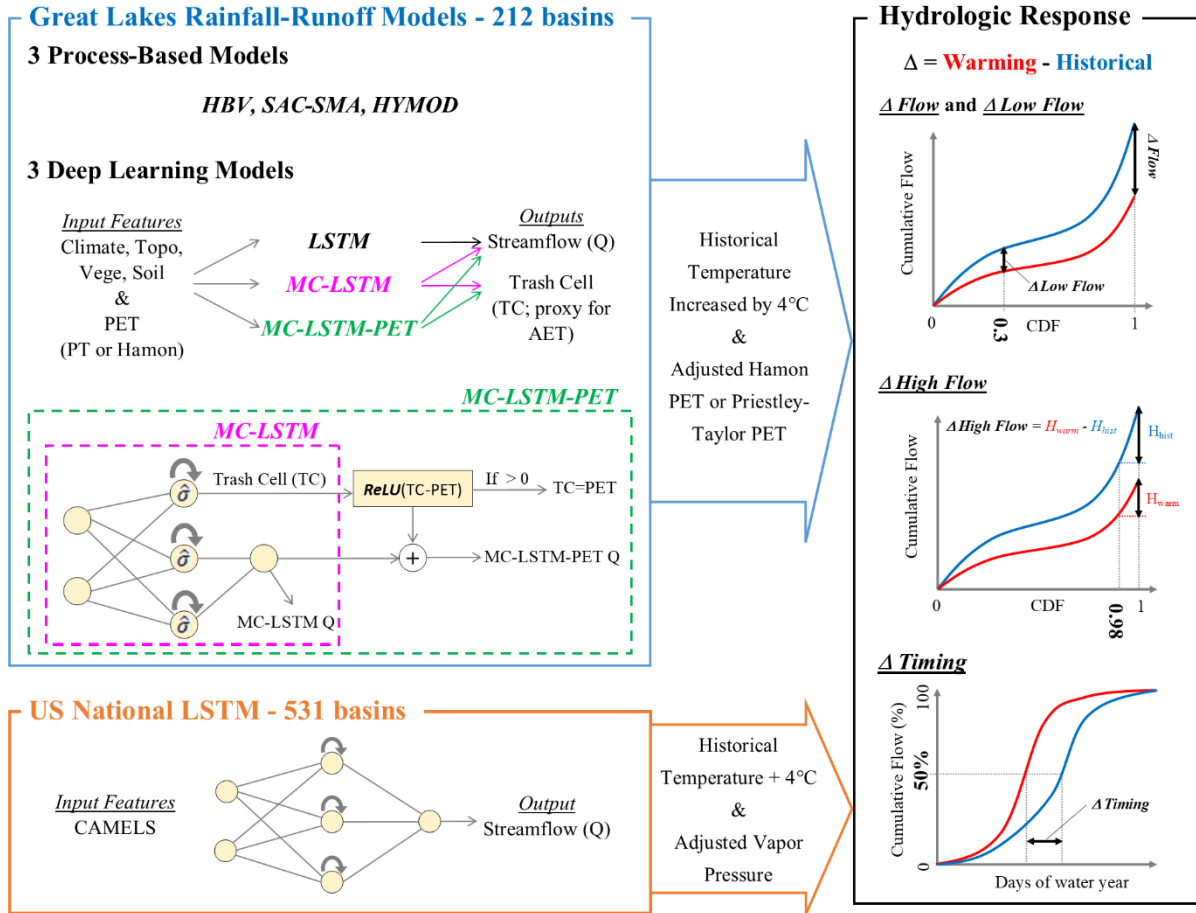
262 rainfall-runoff model estimates of AET.

263

264 **3. Methods**

265 We design an experiment to test the two primary hypotheses of this study, namely that a standard LSTM

266 will overestimate water losses under warming because of an overreliance on historical correlations between

267 temperature and PET, while this effect will be lower in PIML-based rainfall-runoff models designed to

268 better account for water loss in the system. To conduct this experiment, we develop three different DL

269 rainfall-runoff models to predict daily streamflow across the Great Lakes region, as well as three conceptual,

270 process-based models as benchmarks, each of which is trained twice with either an energy budget-based or

271 temperature-based estimate of PET. The DL models include a regional LSTM very similar to the model in

272 Mai et al., (2022), an MC-LSTM that conserves mass, and a new variant of the MC-LSTM that also respects

273   the relationship between PET and water loss (termed MC-LSTM-PET). After comparing historical model

274   performance, we conduct a sensitivity analysis on all models in which historical temperatures are warmed

275   by 4°C, PET is updated based on those warmed temperatures, and all other meteorological variable time

276   series are left unchanged from historical values. This is a similar approach to that taken in SW22, but in

277   contrast to that study this work 1) focuses on the magnitude of streamflow response to warming under two

278   different PET formulations; 2) considers a different set of physics-informed DL models in which the

279   architecture (rather than the inputs or targets) of the model are changed to better preserve physical

280   plausibility under shifts in climate; and 3) evaluates an expanded set of hydrologic metrics to better

281   understand both the plausibility and the variability of responses across the different models. Finally, in a

282   subset of the analysis, we also utilize a fourth DL model, the LSTM used in SW22 that was previously fit

283   to 531 basins across the CONUS (Kratzert et al. 2021), which uses daily precipitation, maximum and

284   minimum temperature, radiation, and vapor pressure as input but not PET. This model is used to evaluate

285   whether a DL model fit to many more watersheds that span a more diverse gradient of climate conditions

286   behaves differently under warming than an LSTM fit only to locations in the Great Lakes basin. Figure 2

287   presents an overview of our experimental design.

288

**Figure 2.** Overview of experiment design. Three deep learning rainfall-runoff models (LSTM, MC-LSTM, MC-LSTM-PET) and three conceptual, process-based models (HBV, SAC-SMA, HYMOD) are trained and tested across 212 watersheds throughout the Great Lakes basin. Models are validated by comparing predictions to streamflow (Q) and actual evapotranspiration (AET). All models are then forced with historical meteorology, but with historical temperatures warmed by 4°C and potential evapotranspiration (PET) updated based on those warmed temperatures using either the Hamon or Priestley-Taylor method. Hydrologic model responses across all models are then compared in terms of long-term mean daily flows, low flows, high flows, and streamflow seasonal timing statistics. The experiment is also repeated with an LSTM fit to 531 basins across the contiguous United States, except that model does not use PET as an input and vapor pressure is also adjusted along with temperature.

## 3.1. Models

### 3.1.1. Benchmark Conceptual Models

We develop three conceptual, process-based hydrologic models as benchmarks, including the Hydrologiska

Byråns Vattenbalansavdelning (HBV) model (Bergström and Forsman, 1973), HYMOD (Boyle, 2001), and

the Sacramento Soil Moisture Accounting (SAC-SMA) model (Burnash, 1995) coupled with SNOW-17

14

(Anderson, 1976). These models are developed as lumped, conceptual models for each watershed, and were selected for several reasons. First, in the Great Lakes Intercomparison Project (Mai et al., 2022), HYMOD was one the best performing process models for both streamflow and AET estimation. SAC-SMA is widely used in the United States, forming the core hydrologic model in NOAA's Hydrologic Ensemble Forecasting System (Demargne et al., 2014). We also found in WS22 that AET from SAC-SMA matched the seasonal pattern of MODIS-derived AET well across California. HBV is also an extremely popular model (Seibert and Bergström, 2022), is used for operational forecasting in multiple countries (Olsson and Lindstrom, 2008; Krøgli et al., 2018), and performs very well in hydrologic model intercomparison projects (Breuer et al., 2009; Plesca et al., 2012; Beck et al., 2016, 2017).

We calibrate the process-based models with the genetic algorithm from Wang et al. (1991) to minimize the mean-squared error (MSE), using a population size equal to 100 times the number of parameters, evolved over 100 generations, and with a spin-up period of 1 year. Each benchmark model is calibrated separately to each of the 141 training sites using the temporal train/test split described in Section 2, and training is repeated 10 separate times with different random initializations to account for uncertainty in the training process and to estimate parametric uncertainty. Benchmark models are developed for the 71 testing sites in two ways: 1) separate models are trained for the testing sites during the training period; and 2) each testing site is assigned a donor from among the 141 training sites, and the calibrated parameters from that donor site are transferred to the testing site. The first of these approaches enables a comparison between DL models fit only to the training sites to benchmark models developed for the testing sites, i.e., a spatial out-of-sample versus in-sample comparison. The second of these approaches enables a more direct spatial out-of-sample comparison between DL and benchmark models. We note that donor sites were used to assign model parameters to testing sites in the benchmarking study of Mai et al. (2022), and to retain direct comparability to the results of that work we use the same donor sites for each testing site. Donor sites were selected based on spatial proximity, while also prioritizing donor sites that were nested within the watershed of the testing site.

332

### 3.1.2. LSTM

334    We develop a single, regional LSTM for predicting daily streamflow across the Great Lakes region. In the

335    LSTM, nodes within hidden layers feature gates and cell states that address the vanishing gradient problem

336    of classic recurrent neural networks and help capture long-term dependencies between input and output

337    time series. The model defines a $D$-dimensional vector of recurrent cell states $\boldsymbol{c}[t]$ that is updated over a

338    sequence of $t$=1,…,$T$ time steps based on a sequence of inputs $\boldsymbol{x} = \boldsymbol{x}[1], \dots, \boldsymbol{x}[T]$, where each input $\boldsymbol{x}[t]$ is

339    a $K$-dimensional vector of features. Information stored in the cell states is then used to update a $D$-

340    dimensional vector of hidden states $\boldsymbol{h}[t]$, which form the output of the hidden layer in the model. The

341    structure of the LSTM is given as follows:

342

343    $\boldsymbol{i}[t] = \sigma(\boldsymbol{W}_i\boldsymbol{x}[t] + \boldsymbol{U}_i\boldsymbol{h}[t-1] + \boldsymbol{b}_i)$         (Eq. 1.1)

344    $\boldsymbol{f}[t] = \sigma\big(\boldsymbol{W}_f\boldsymbol{x}[t] + \boldsymbol{U}_f\boldsymbol{h}[t-1] + \boldsymbol{b}_f\big)$         (Eq. 1.2)

345    $\boldsymbol{g}[t] = tanh\big(\boldsymbol{W}_g\boldsymbol{x}[t] + \boldsymbol{U}_g\boldsymbol{h}[t-1] + \boldsymbol{b}_g\big)$         (Eq. 1.3)

346    $\boldsymbol{o}[t] = \sigma(\boldsymbol{W}_o\boldsymbol{x}[t] + \boldsymbol{U}_o\boldsymbol{h}[t-1] + \boldsymbol{b}_o)$         (Eq. 1.4)

347    $\boldsymbol{c}[t] = \boldsymbol{f}[t] \odot \boldsymbol{c}[t-1] + \boldsymbol{i}[t] \odot \boldsymbol{g}[t]$         (Eq. 1.5)

348    $\boldsymbol{h}[t] = \boldsymbol{o}[t] \odot tanh(\boldsymbol{c}[t])$         (Eq. 1.6)

349    $\boldsymbol{y}[T] = ReLU\big(\boldsymbol{W}_y\boldsymbol{h}[T] + b_y\big)$         (Eq. 1.7)

350

351    Here, the input gate ($\boldsymbol{i}[t]$) controls how candidate information ($\boldsymbol{g}[t]$) from inputs and previous hidden states

352    flows to the current cell state ($\boldsymbol{c}[t]$); the forget gate ($\boldsymbol{f}[t]$) enables removal of information within the cell

353    state over time; and the output gate ($\boldsymbol{o}[t]$) controls information flow from the current cell state to the hidden

354    layer output. All bolded terms are vectors, and $\odot$ denotes element-wise multiplication. To produce

355    streamflow predictions, $\boldsymbol{h}[T]$ at the last time step in the sequence is passed through a fully connected layer

356    to a single-node output layer (i.e., a many-to-one formulation). We ensure nonnegative streamflow

357    predictions using the rectified linear unit (ReLU) activation function for the output neuron, expressed as

358    ReLU(*x*) = max(0,*x*). Importantly, there are no constraints requiring the mass of water entering as

359    precipitation to be conserved within this architecture.

360

361    The LSTM takes *K*=39 input features: 9 dynamic and 30 static. The dynamic input features are basin-

362    averaged climate, including daily precipitation, maximum temperature, minimum temperature, net

363    incoming shortwave radiation, specific humidity, surface air pressure, zonal and meridional components of

364    wind, and PET. The static features represent catchment attributes (see Table 1) and are repeated for all time

365    steps in the input sequences $\boldsymbol{x}$. All input features are standardized before training (by subtracting the mean

366    and dividing by the standard deviation for data across all training sites in the training period). Note that we

367    do not standardize the observed streamflow, besides dividing by drainage area to represent streamflow in

368    units of millimeters.

369

370    We train the LSTM by minimizing the mean-squared error averaged over the 141 training watersheds

371    during the training period:

$$MSE = \frac{1}{N}\sum_{n=1}^{N}\frac{1}{T_n}\sum_{t=1}^{T_n}\left(\hat{Q}_{n,t} - Q_{n,t}\right)^2 \tag{2}$$

373    where $N$ is the number of training watersheds and $T_n$ is the number samples in the $n^{th}$ watershed. $\hat{Q}_{n,t}$ and

374    $Q_{n,t}$ are, respectively, the streamflow prediction and observation for basin $n$ and day $t$. To estimate $\hat{Q}_{n,t}$,

375    we feed into the network an input sequence for the past $T$=365 days. The model was developed with 1

376    hidden layer composed of $D$=256 nodes, a mini-batch size of 256, a learning rate of 0.0005, and a drop-out

377    rate of 0.4, and it was trained across 30 epochs. All hyperparameters (number of hidden layer nodes, mini-

378    batch size, learning rate, dropout rate, and number of epochs) were selected in a 5-fold cross-validation on

379    the training sites (see Table S2 for details on grid search). Network weights are tuned using the ADAM

380    optimizer (Kingma & Ba, 2015). The model is trained 10 separate times with different random

381    initializations to account for uncertainty in the training process.

382

383    For the evaluation of streamflow responses to warming, we also use an LSTM taken from Kratzert et al.

384    (2021) and employed in SW22, which was fit to 531 basins across the contiguous United States (hereafter

385    called the National LSTM). This model was trained using a different set of data compared to our Great

386    Lakes LSTM but also used a mix of dynamic and static features, all of which were drawn from the

387    Catchment Attributes and Meteorology for Large-Sample Studies (CAMELS) dataset (Newman et al.,

388    2015). This model uses daily precipitation, maximum and minimum temperature, shortwave downward

389    radiation, and vapor pressure as input but not PET. However, we note that temperature, radiation, and vapor

390    pressure are the three major inputs (besides wind speeds) needed to calculate energy budget-based PET.

391    There are 29 CAMELS watersheds located within the Great Lakes basin, and 17 of those 29 watersheds

392    were also used in the training and testing sets for the Great Lakes LSTM (see Figure 1).

393

### 394    3.1.3. MC-LSTM

395    Following Hoedt et al. (2021) and Frame et al. (2021b), we adapt the architecture of the LSTM into a mass

396    conserving MC-LSTM that preserves the water balance within the model, i.e., the total quantity of

397    precipitation entering the model is tracked and redistributed to streamflow and losses from the watershed.

398    Using similar notation as for the LSTM above, the model structure is given as follows:

399

400    $\hat{c}[t-1] = \frac{c[t-1]}{\|c[t-1]\|_1}$     (Eq. 3.1)

401    $i[t] = \hat{\sigma}(W_i x[t] + U_i \hat{c}[t-1] + V_i a[t] + b_i)$     (Eq. 3.2)

402    $o[t] = \sigma(W_o x[t] + U_o \hat{c}[t-1] + V_o a[t] + b_o)$     (Eq. 3.3)

403    $R[t] = \hat{\sigma}(W_R x[t] + U_R \hat{c}[t-1] + V_R a[t] + b_R)$     (Eq. 3.4)

404    $m[t] = R[t]c[t-1] + i[t]x[t]$     (Eq. 3.5)

405    $c[t] = (1 - o[t]) \odot m[t]$     (Eq. 3.6)

406    $h[t] = o[t] \odot m[t]$     (Eq. 3.7)

407

408 Here, the inputs to the model are split between quantities $\mathbf{x}[t]$ to be conserved (i.e., precipitation), and non-

409 conservative inputs $\mathbf{a}[t]$ (i.e., temperature, wind speeds, PET, catchment properties, etc.). Water in the

410 system is stored in the $D$-dimensional vector $\mathbf{m}[t]$ and is updated at each time step based on water left over

411 from the previous time step ($\mathbf{c}[t\text{-}1]$) and water entering the system at the current time step ($\mathbf{x}[t]$). The input

412 gate $\mathbf{i}[t]$ and a redistribution matrix $\mathbf{R}[t]$ are designed to ensure water is conserved from $\boldsymbol{c}[t-1]$ and $\boldsymbol{x}[t]$

413 to $\mathbf{m}[t]$, by basing these quantities on a normalized sigmoid activation function:

414

415 $\quad \hat{\sigma}(z_j) = \frac{\sigma(z_j)}{\sum_j \sigma(z_j)}$ (Eq. 4)

416

417 Here, $\sigma(\cdot)$ is the sigmoid activation function, while $\hat{\sigma}(\cdot)$ is a normalized sigmoid activation that produces a

418 vector of fractions that sum to unity.

419

420 The mass in $\boldsymbol{m}[t]$, which is stored across $D$ elements in the vector, is then distributed to the output of the

421 hidden layer, $\boldsymbol{h}[t]$, or the next cell state, $\boldsymbol{c}[t]$. To account for water losses from evapotranspiration or other

422 sinks, one element of the $D$-dimensional vector $\boldsymbol{h}[t]$ is considered a 'trash cell', and the output of this cell

423 is ignored when calculating the final streamflow prediction, which at time $T$ is given by the sum of outgoing

424 water mass:

425

426 $\quad y[T] = \sum_{d=1}^{D-1} h_d[T]$ (Eq. 5)

427

428 Here, the $D^{\text{th}}$ cell of $\boldsymbol{h}$ ($h_D$) is set as the trash cell, and water allocated to this cell at each time step $t=1,..,T$

429 is lost from the system. We note that the MC-LSTM was trained in the same way as the LSTM (i.e., same

430 inputs, loss function, training and test sets, hyperparameter selection process, number of ensemble members

431 with random initialization).

432

### 3.1.4. MC-LSTM-PET

434    We also propose a novel variant of the MC-LSTM that requires water lost from the system to not exceed

435    PET (hereafter referred to as the MC-LSTM-PET). In the original MC-LSTM, any amount of water can be

436    delegated to the trash cell $h_D$. Therefore, while water is conserved in the MC-LSTM, the model has the

437    freedom to transfer any amount of water from $\boldsymbol{m}[t]$ to the trash cell (and out of the hydrologic system) as

438    it seeks to improve the loss function during training. This has the benefit of handling biased data, e.g., cases

439    where the precipitation input to the system is systematically too high compared to the measured outflow.

440    However, this structure also has the drawback of potentially removing more water from the system than is

441    physically plausible. To address this issue, we propose a small change to the architecture of the MC-LSTM,

442    where any water relegated to the trash cell that exceeds PET at time $t$ is directed back to the stream:

443

444    $$y[t] = \sum_{d=1}^{D-1} h_d[t] + ReLU(h_D[t] - PET[t]) \qquad \text{(Eq. 6)}$$

445

446    Here, the ReLU activation ensures that any water in the trash cell ($h_D$) which exceeds PET at time $t$ is

447    added to the streamflow prediction $y[t]$, but the streamflow prediction is the same as the original MC-

448    LSTM (Eq. 5) if water in the trash cell is less than PET. This approach assumes that the maximum allowable

449    water lost from the system cannot exceed PET, and therefore ignores other potential terminal sinks (e.g.,

450    inter-basin lateral groundwater flows; human diversions and inter-basin transfers). This assumption is more

451    strongly supported in moderately-sized ($> 200$ km$^2$), low-gradient, non-arid watersheds where inter-basin

452    groundwater flows are less impactful (Fan 2019; Gordon et al., 2022), such as the Great Lakes basins

453    examined in this work. However, we discuss the potential to relax the assumptions of the MC-LSTM-PET

454    model in Section 5. The MC-LSTM-PET was trained in the same way as the LSTM (i.e., same inputs, loss

455    function, training and test sets, hyperparameter selection process, number of ensemble members with

456    random initialization).

457

## 3.2. Model Performance Evaluation

As noted previously, 141 of the watersheds are designated as training sites, and the remaining 71 watersheds are used for testing. In addition, the training and testing periods were restricted to January 2000 -December 2010 and January 2011 – December 2017, respectively. This provides three separate ways to evaluate model performance:

- Temporal validation - Performance across models is evaluated at training sites during the testing period.

- Spatial validation - Performance across models is evaluated at testing sites during the training period.

- Spatiotemporal validation - Performance across models is evaluated at testing sites during the testing period.

All three evaluation strategies are utilized. For benchmark process-based models that are calibrated locally on a site-by-site basis, we consider model versions that are transferred to testing sites from training sites, as well as models that are trained to the testing sites directly (see Section 3.1.1). The former can be used for all three evaluation strategies above, while the latter can only be used for temporal validation at the testing sites.

Following other intercomparison studies (Frame et al., 2022; Gauch et al., 2021a; Klotz et al., 2022; Kratzert et al., 2021), several metrics are considered for model evaluation, including percent bias (PBIAS), the Nash-Sutcliffe Efficiency (NSE; Nash and Sutcliffe, 1970), Kling-Gupta Efficient (KGE; Gupta et al. 2009), top 2% peak flow bias (FHV; Yilmaz et al. 2008), and bottom 30% low flow bias (FLV; Yilmaz et al. 2008). Each metric is calculated separately for training and testing periods for each site. For all models, all results are estimated from the ensemble mean from 10 separate training trials.

482

483    For the process models, the MC-LSTM, and the MC-LSTM-PET, we also compare simulations of AET to

484    AET from the GLEAM database. We note that AET data were not used to train any of the models. For the

485    process models, AET is a direct output of the model and so can immediately be extracted for comparison,

486    but AET is not directly simulated by the MC-LSTM or MC-LSTM-PET. Instead, we assume water

487    delegated to the trash cell permanently leaves the system because of evapotranspiration. Several metrics

488    are used to compare model based AET to GLEAM AET, including KGE, correlation, and PBIAS, and the

489    comparison is conducted for training sites during the training period and under temporal, spatial, and

490    spatiotemporal validation (as described above). Similar to streamflow, all AET results are based on the

491    ensemble mean from the 10 separate training trials.

492

493    **3.3. Evaluating Hydrologic Response under Warming**

494    All Great Lakes models in this study are trained twice with different PET estimates as input, including the

495    Hamon method (a temperature-based approach; Hamon, 1963) and the Priestley-Taylor method (an energy

496    budget-based approach; Priestley and Taylor, 1972). We select the Hamon method because of its stronger

497    dependence on temperature compared to other temperature-based approaches that also depend on radiation

498    (e.g., Hargreaves and Samani, 1985; Oudin et al., 2005). We select the Priestley-Taylor method based on

499    its widespread use in the literature (Wu et al., 2021; Su and Singh, 2023) and its approximation of the more

500    physically-based Penman-Monteith approach (Allen et al. 1998). Together, these two approaches lie

501    towards the lower and upper bounds of temperature sensitivity across multiple PET approaches (see Shaw

502    and Riha, 2011).

503

504    PET (in mm/day) under the Hamon method is calculated as follows (Shaw and Riha, 2011):

505

506    $PET_H = \alpha_H \times 29.8 \times Hr_{day} \frac{e_{sat}}{T_a + 273.2}$                              (Eq. 7)

507 $$e_{sat} = 0.611 \times exp\left(\frac{17.27 \times T_a}{237.3 + T_a}\right) \qquad \text{(Eq. 8)}$$

508  where $Hr_{day}$ is the number of daylight hours, $T_a$ is the average daily temperature (°C) calculated from

509  daily minimum and maximum temperature, $e_{sat}$ is the saturation vapor pressure (kPa), and $\alpha_H$ is a

510  calibration coefficient set to 1.2 for all models in this study (similar to Lu et al., 2005).

511

512  PET under the Priestley-Taylor method is calculated as follows:

513

514 $$PET_{PT} = \alpha_{PT}\left(\frac{\Delta(T_a) \times (R_n - G)}{\lambda(\Delta(T_a) + \gamma)}\right) \times 1000 \qquad \text{(Eq. 9)}$$

515

516  Here, $\Delta(T_a)$ is the slope of the saturation vapor pressure temperature curve (kPa/°C) and is a function of

517  $T_a$, $\gamma$ is the psychrometric constant (kPa/°C), $\lambda$ is the volumetric latent heat of vaporization (MJ/m³), $R_n$ is

518  the net radiation (MJ/m²-day) equal to the difference between net incoming shortwave ($R_{ns}$) and net

519  outgoing longwave ($R_{nl}$) radiation, $G$ is the heat flux to the ground (MJ/m²-day), and $\alpha_{PT}$ is a dimensionless

520  coefficient set to 1.1 for all models in this study (similar to Szilagyi et al., 2017). Details on how to calculate

521  $\gamma$, $\Delta(T_a)$, and $R_{nl}$ are available in Allen et al. (1998), and we assume $G$=0. Net shortwave radiation is given

522  by $R_{ns} = (1 - \zeta)R_s$, with $\zeta = .23$ the assumed albedo and $R_s$ the incoming shorwave radiation. We note

523  that net outgoing longwave radiation $R_{nl}$ is a function of maximum and minimum temperature, actual vapor

524  pressure, and $R_s$ (see Eq. 39 in Allen et al. 1998). All exogenous meteorological inputs for the two methods

525  are derived from the Regional Deterministic Reanalysis System v2 (see Section 2). We note that using

526  $\alpha_H = 1.2$ and $\alpha_{PT} = 1.1$ leads to very similar long-term average PET estimates between the Hamon and

527  Priestley-Taylor methods under baseline climate conditions, helping to ensure their comparability. We also

528  note that both PET series are highly correlated with daily average temperatures (average Pearson

529  correlations across sites of 0.94 and 0.83 for Hamon and Priestley-Taylor PET, respectively).

530

531  We then conduct a sensitivity analysis of model response in which the historical minimum and maximum

532  temperature time series are increased uniformly by 4 °C, and the two PET estimates are updated using these

533  warmed temperatures. We focus the assessment on training period data at the training sites, so that any

534  differences in responses that emerge between the DL and process models are due to model structural

535  differences and not the effects of spatiotemporal regionalization. In the Priestly-Taylor method, we maintain

536  historical values for $R_s$ to isolate how changes in temperature and its effect on $\Delta(T_a)$ and $R_{nl}$ influence

537  changes in PET. The use of historical $R_s$ is supported by the results from CMIP5 projections presented in

538  Lai et al. (2022), but this assumption is discussed further in Section 5.

539

540  We also conduct a similar sensitivity analysis on the National LSTM, which uses five dynamic input

541  features from the CAMELS dataset (daily precipitation, maximum temperature, minimum temperature, $R_s$,

542  and water vapor pressure). Here, temperatures are warmed by 4°C, while precipitation and $R_s$ are held at

543  historical values. There is a strong correlation between vapor pressure and minimum temperature in the

544  CAMELS dataset, since minimum temperature is used to estimate the water vapor pressure (Newman et al.,

545  2015). Thus, to run the National LSTM under warming, we also adjust the vapor pressure input based on

546  the change imposed to minimum temperature. This procedure is detailed in SW22.

547

548  For both the Great Lakes DL models and the National LSTM, the dynamic inputs are adjusted based on the

549  warming scenarios above. We also consider changes to the static input features that depend on temperature

550  and PET in their calculation (e.g., pet_mean, aridity, t_mean, frac_snow; see Table 1 for feature descriptions

551  and Supporting Information S1 and Table S1 for details on adjustments to these features), and then run all

552  models using two settings: 1) with changes only to the dynamic features, and 2) with changes to both

553  dynamic features and to static features that depend on those dynamic features. In total, there are six

554  scenarios run in this work, which are shown in Table 2.

555

556 **Table 2**. Overview of the setup for the different scenarios run in this analysis. All models are driven with
557 temperatures warmed by 4°C. The Great Lakes models include the HBV, SAC-SMA, HYMOD, LSTM,
558 MC-LSTM, and MC-LSTM models that are trained and tested to the 212 sites across the Great Lakes basin.
559

| Scenario | Model | PET method adjusted with warmer temperatures | Are static features also changed along with dynamic features? |
|---|---|---|---|
| 1 | Great Lakes models | Hamon | Yes |
| 2 | Great Lakes models | Priestley-Taylor | Yes |
| 3 | Great Lakes models | Hamon | No |
| 4 | Great Lakes models | Priestley-Taylor | No |
| 5 | National LSTM | NA | Yes |
| 6 | National LSTM | NA | No |

560

561

562 Ultimately, for each model we compare hydrologic responses under the warmed scenario to their values

563 under the baseline scenario with no warming. For the National LSTM, we only consider basins in the

564 CAMELS dataset within the Great Lakes Basin. For the process models, we also evaluate the uncertainty

565 in hydrologic response based on the range predicted across the 10 different training trials, as a simple means

566 to evaluate how parametric uncertainty influences the predictions. We examine four different metrics for

567 this comparison, including:

568  • AVG.Q: the long-term mean of daily streamflow across the entire series.

569  • FHV: the average of the top 2% peak flows.

570  • FLV: the average of the bottom 30% low flows.

571  • COM: the median center of mass across all water years, where the center of mass is defined as the

572   day of the water year by which half of the total annual flow has passed.

573

574 If our hypothesis is correct that the LSTM cannot distinguish water loss differences with different PET

575 series but similar warming while process-based and PIML models can, we would expect that under the

576 LSTM using both PET series, long-term mean flow will decline substantially and with similar magnitude

577 to the process models using the temperature-based PET method but not the energy budget-based PET

578 method. We would also expect the National LSTM to exhibit similar behavior, even though it was able to

579 learn from a larger set of watersheds across a more diverse range of climate conditions. Finally, if our

25

580  hypothesis is correct, we would expect the PIML models (MC-LSTM, MC-LSTM-PET) to follow the

581  process model responses more closely across the two different PET series, at least in terms of the difference

582  in magnitude of long-term mean streamflow declines. To facilitate a broader inter-model comparison of DL

583  and process-based models under warming (which is largely absent from the literature), we also explore the

584  differences in low flow (FLV), high flow (FHV), and seasonal timing (COM) metrics across all model

585  versions, where we have less reason to anticipate how DL and process models will differ in their responses

586  and across PET formulations. However, for responses like seasonal streamflow timing (COM), we do

587  anticipate that realistic responses should show a shift towards more streamflow earlier in the year, as

588  warmer temperatures lead to more precipitation falling as rain rather than snow and drive snowmelt earlier

589  in the spring.

590

591  **4. Results**

592  **4.1. Model Performance Evaluation**

593  Figure 3 shows the distribution of KGE values across sites for streamflow from the LSTM, MC-LSTM,

594  MC-LSTM-PET, and the three process-based models for both the training and testing sites during both the

595  training and testing periods. All results here and elsewhere in Section 4.1 are shown for the models fit with

596  Priestley-Taylor PET, but there is little difference in performance for the models fit with Hamon PET (see

597  Figure S1). For the process-based models, we show results for models fit to the training sites and then used

598  as donors at the testing sites, as well as models fit to the testing sites directly. We denote the latter with the

599  suffix "-test" and note that performance metrics at the training sites are not available for process models fit

600  to the testing sites.

601

602  Several insights emerge from Figure 3. First, for the training sites during the training period, all models

603  perform very well (Figure 3a). Across the three process models, the median KGE is 0.79, 0.78, and 0.77

604  for HBV, SAC-SMA, and HYMOD, respectfully. However, unsurprisingly, the DL models perform better

26

605　for the training data, with median KGE values all equal or above 0.88. The LSTM performs best in this

606　case. Under temporal validation (training sites during the testing period), performance degrades somewhat

607　across all models, and the differences in KGE between all process-based models and between all DL models

608　shrink considerably (Figure 3c). Larger performance declines are seen at the testing sites during the training

609　period (Figure 3b) and testing period (Figure 3d). Here, the median KGE for all process models falls to

610　between 0.54-0.58 when streamflow at the testing sites is estimated with donor models from nearby gauged

611　watersheds. In contrast, process models fit to the testing sites (denoted "-test") exhibit performance similar

612　to that seen in Figure 3a,c. All three DL models perform quite well for the testing sites, with median KGE

613　values above 0.71 in both time periods. This is only modestly below the median KGE for the process models

614　fit to the testing sites, which is quite impressive given that this represents the spatial out-of-sample

615　performance of the DL models. We even see that for approximately 20% of testing sites during the training

616　period, the DL models outperform the process models fit to those locations in that period.

617

**Figure 3.** The distribution of Kling-Gupta efficiency (KGE) for streamflow estimates across sites from each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results for the testing period are shown in panels (c) and (d), respectively. For the process models fit to the testing sites (denoted "-test"), no performance results are available at the training sites. All models are trained using Priestley-Taylor PET.

Table 3 shows the median KGE, NSE, PBIAS, FHV, and FHL across testing sites for all models, excluding the process models fit to the testing sites. Similar to Figure 3, all three DL models outperform the donor-based process models at the testing sites for all metrics. The performance across the three different DL models is similar, although there are some notable differences. In particular, the LSTM outperforms the MC-LSTM and MC-LSTM-PET for NSE and FLV (as well as KGE in the training period), the MC-LSTM-PET outperforms the LSTM and MC-LSTM for PBIAS, and either the MC-LSTM or MC-LSTM-PET are the best performers for FHV. The fact that the MC-LSTM-PET performs best for PBIAS of all models suggests that the PET constraint imposed in that model improves the overall accounting of water entering
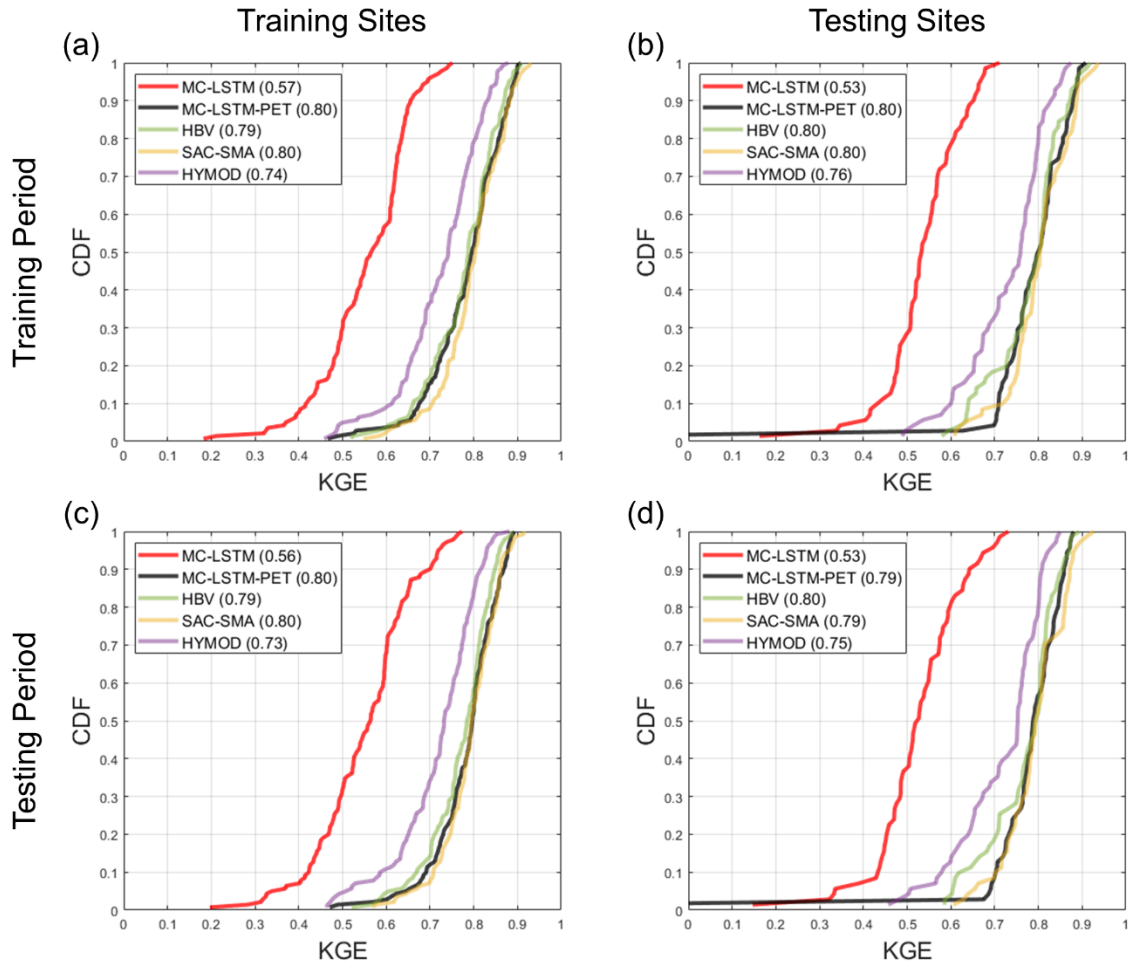
633  and existing the watershed on a long-term basis. We also note that percent biases for FLV are high because

634  the absolute magnitude of low flows is small, so small absolute biases still lead to large percent biases.

635

636  **Table 3.** The median KGE, NSE, PBIAS, FHV, and FLV for streamflow across testing sites for the training
637  and testing periods for all models (excluding the process models fit to the testing sites). The metric from
638  the best performing model in each period is bolded. All models are trained using Priestley-Taylor PET.

| Model | Testing Sites: Training Period | | | | | Testing Sites: Testing Period | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | KGE | NSE | PBIAS | FHV | FLV | KGE | NSE | PBIAS | FHV | FLV |
| LSTM | **0.76** | **0.77** | 9.66 | 17.58 | **30.98** | **0.72** | **0.68** | 12.15 | 26.01 | **27.32** |
| MC-LSTM | 0.74 | 0.72 | 9.48 | **15.52** | 41.46 | **0.72** | 0.65 | 12.13 | 22.82 | 35.80 |
| MC-LSTM-PET | 0.73 | 0.72 | **8.63** | 18.80 | 48.10 | 0.71 | 0.66 | **10.22** | **22.49** | 44.43 |
| HBV | 0.58 | 0.50 | 9.99 | 32.22 | 63.96 | 0.55 | 0.50 | 12.68 | 34.76 | 57.20 |
| SAC-SMA | 0.57 | 0.48 | 11.74 | 34.72 | 45.17 | 0.54 | 0.47 | 12.24 | 40.45 | 46.78 |
| HYMOD | 0.58 | 0.48 | 10.07 | 33.68 | 58.06 | 0.54 | 0.48 | 12.52 | 36.07 | 60.32 |

639

640  Figure 4 shows similar results as Figure 3, but for the KGE based on estimates of AET. Also, only donor

641  process models are shown for the testing sites. Results for correlation and PBIAS are available in the

642  Supplemental Information (Figures S2-S3). Here, the LSTM is not included because estimates of AET are

643  unavailable, while AET from the MC-LSTM and MC-LSTM-PET is based on water relegated to the trash

644  cell. Note that none of the models were trained for AET, and so results at training sites during the training

645  period also provide a form of model validation. Figure 4 shows that SAC-SMA and HBV predict AET with

646  relatively high degrees of accuracy for both training and testing sites in both periods (median KGE between

647  0.79-0.80). Performance is slightly worse for HYMOD. Notably, the MC-LSTM-PET exhibits very similar,

648  strong performance for all sites and periods as compared to SAC-SMA and HBV, except for one testing

649  site. In contrast, the MC-LSTM performs the worst of all models, with median KGE values ranging between
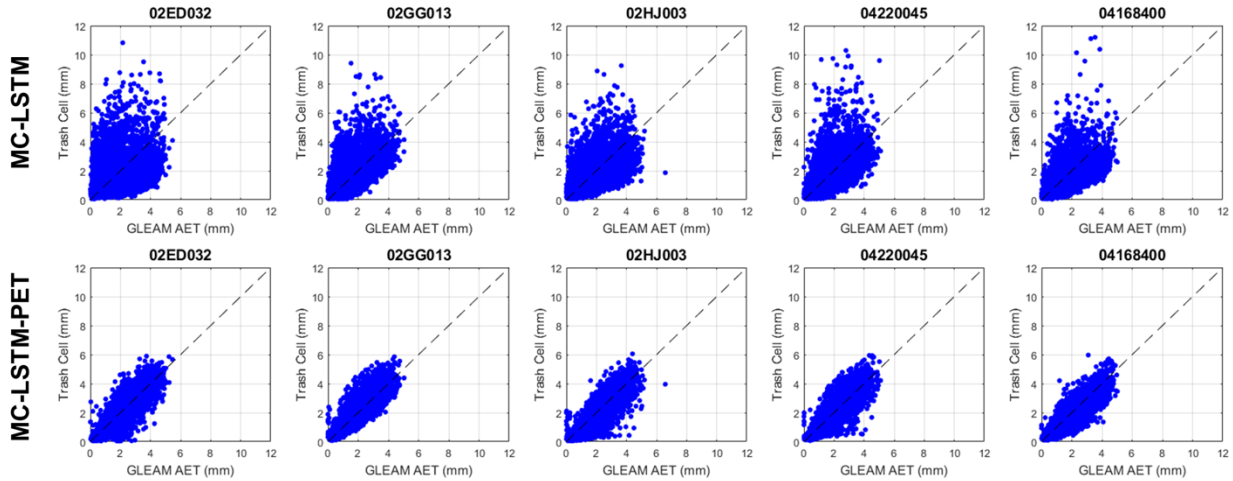
650  0.53-0.57.

651

**Figure 4.** The Kling-Gupta efficiency (KGE) for AET estimated from each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results for the testing period are shown in panels (c) and (d), respectively. The LSTM is not included in this comparison. All models are trained using Priestley-Taylor PET.

Further investigation reveals that the differences in KGE between the MC-LSTM and MC-LSTM-PET models for AET are largely driven by differences in correlation (see Figure S2). We examine this difference in more detail in Figure 5, which presents scatterplots of GLEAM AET versus water allocations to the trash cell for the two models from five randomly sampled testing sites across both training and testing periods (see Figure 1; also Table S3). Trash cell water from the MC-LSTM is not only more scattered around GLEAM AET compared to the MC-LSTM-PET, but it also exhibits many outlier values that are two to five times larger than GLEAM AET. The MC-LSTM-PET follows the variability of GLEAM AET much more closely, with virtually no outliers that exceed GLEAM AET by large margins. This suggests that the PET

30

666    constraint on the trash cell in the MC-LSTM-PET helps water allocated to that cell more faithfully represent

667    evaporative water loss in the DL model.
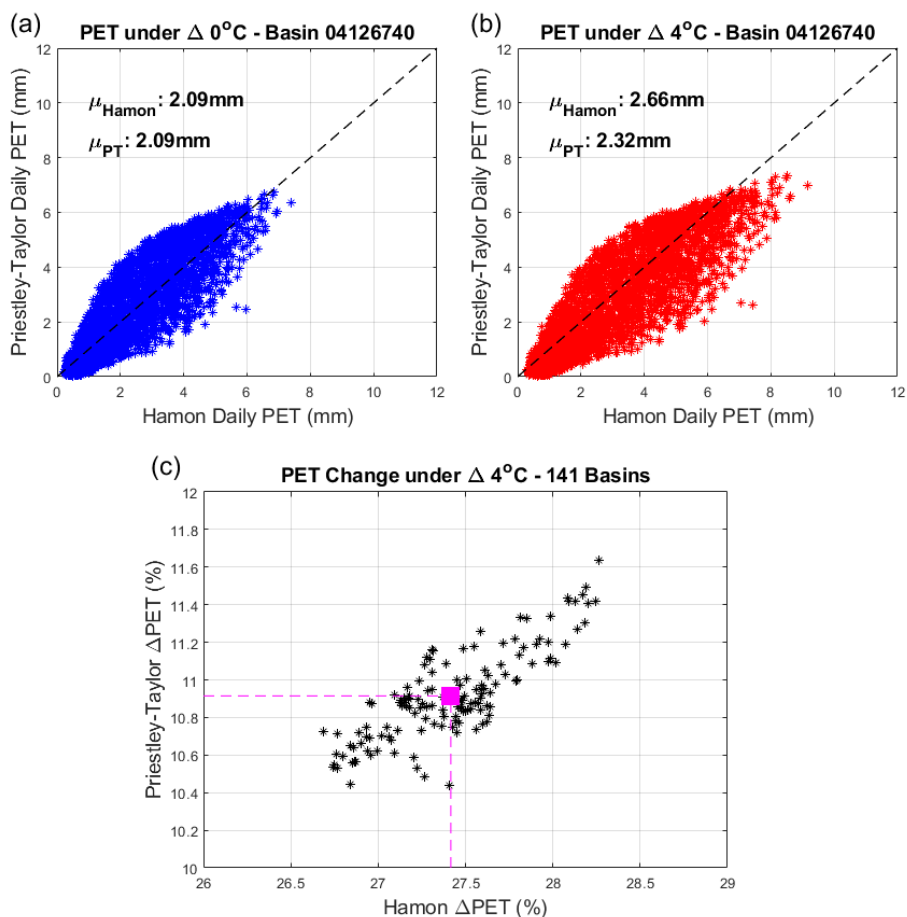
668



669

670    **Figure 5.** Scatterplots of daily AET versus trash cell water for the (top) MC-LSTM and (bottom) MC-
671    LSTM-PET at five randomly selected testing sites across both training and testing periods. All models are
672                            trained using Priestley-Taylor PET.

673

## 4.2. Evaluating Hydrologic Response under Warming

675    Next, we evaluate streamflow responses under a 4 °C warming scenario. We focus on training sites during

676    the training period, so that any differences that emerge between DL and process models are only related to

677    model structure and not spatiotemporal regionalization. However, our results are largely unchanged if based

678    on responses for testing sites in the testing period (see Figure S4). First, we show the differences in historic

679    and warming-adjusted PET when using the Hamon and Priestley-Taylor methods (Figure 6). For the

680    training period without any temperature change, PET estimated from the two methods is very similar

681    (Figure 6a; shown at one sample location for demonstration, see Figure 1 and Table S3). However, under

682    the scenario with 4 °C of warming, Hamon-based PET is substantially larger than Priestley-Taylor based

683    PET (Figure 6b). On average, this difference reaches ~16% across all training sites and exhibits very little

684    variability across locations (Figure 6c). The primary reason for the difference in the estimated change in

685    PET is that the Hamon method attributes PET entirely to temperature, while only a portion of PET is based

686     on temperature in the Priestley-Taylor method, with the rest based on $R_n$. It is worthwhile to note that $R_n$

687     does increase with temperature through its effects on net outgoing longwave radiation, but these changes

688     are generally less than 5% across all sites (Allen et al. 1998).

689



690

691        **Figure 6.** (a) Daily PET estimated using the Hamon and Priestley-Taylor method for one sample
692 watershed, under historic climate conditions in the training period. (b) Same as (a), but under the scenario
693 with 4 °C of warming. (c) Percent change in average PET with 4 °C of warming across all training sites
694 using the Hamon and Priestley-Taylor methods.
695

696 Figure 7 shows how these differences in PET under warming propagate into changes in different attributes

697 of streamflow across training sites in the training period. The left and right columns of Figure 7 show

698 streamflow responses using Hamon and Priestley-Taylor PET, respectively, while the rows of Figure 7

699   show the distribution of changes in different streamflow attributes (AVG.Q, FLV, FHV, COM) across

700   models. Figure 7 shows results for DL models where only the dynamic inputs are changed under warming.

701

702   Starting with changes in AVG.Q, Figure 7a,b shows that under the Hamon method for PET, the DL models

703   exhibit similar changes in long-term mean streamflow to the process-based models, with the median

704   ΔAVG.Q across sites ranging between -17% and -25% across all models. However, when using Priestley-

705   Taylor PET, larger differences in the distribution of ΔAVG.Q emerge. Across all three process models, the

706   median ΔAVG.Q is between -6% to -9%, and very few locations exhibit ΔAVG.Q less than -20%.

707   Conversely, the LSTM shows a median water loss of -20% under Priestley-Taylor PET and a very similar

708   distribution of water losses regardless of whether Hamon or Priestley-Taylor PET was used. The MC-

709   LSTM is also relatively insensitive to PET, and as compared to the process models, the MC-LSTM tends

710   to predict smaller absolute changes to AVG.Q for Hamon PET and larger changes under Priestley-Taylor

711   PET. Only the MC-LSTM-PET model achieves water loss that is considerably smaller under Priestley-

712   Taylor PET than Hamon PET and closely follows the process models in both cases.

713

714   The overall pattern of change in low flows (FLV) is very similar across all three DL models, with median

715   declines between -15% to -25% and little variability across sites (Figure 7c,d). The process models disagree

716   on the sign of change for FLV, and also bound the changes predicted by the DL models. HBV and HYMOD

717   show mostly increases to FLV under warming and Priestley-Taylor PET, and a mix of increases and

718   decreases across sites for Hamon PET. SAC-SMA exhibits large declines in FLV under warming and

719   Hamon PET, and shows a median change that is similar to the DL models under Priestley-Taylor PET. The

720   percent changes in FLV across models tend to be large because the absolute magnitude of FLV is small,

721   and so small changes in millimeters of flow lead to large percent changes. This can be seen in sample daily

722   hydrographs for two sites (see Figure S5), where visually the changes in low flows are difficult to discern
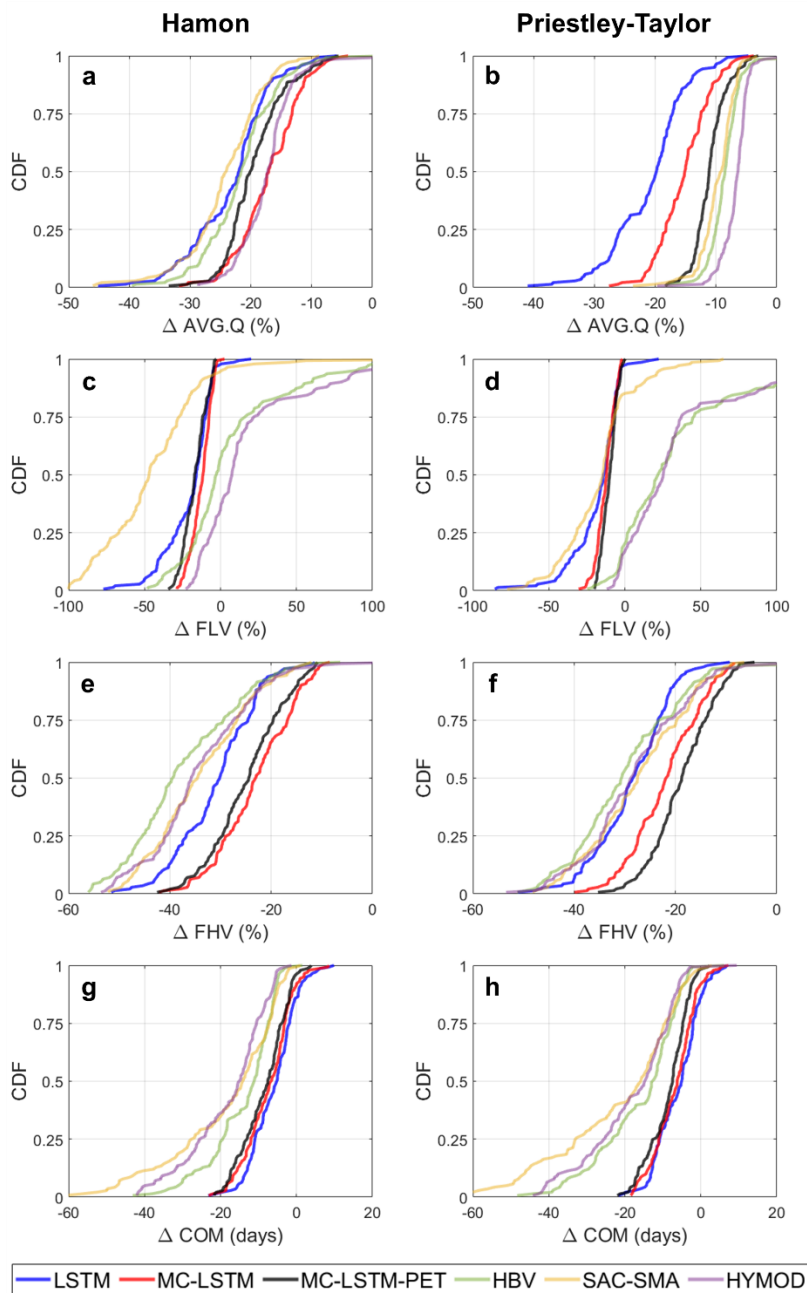
723   because they are all near zero for all models, but the change in the FLV statistic varies significantly across

724   the six models and two sites (-56% to +40%).

725

726   The differences between process-based and DL simulated changes for high flows (FHV; Figure 7e,f) and

727   seasonal timing (COM; Figure 7g,h) are relatively consistent, with the process models exhibiting more

728   substantial declines in high flows and earlier shifts in seasonal timing compared to the DL models. The

729   choice of PET method has an impact on process-model based changes in FHV, with larger declines under

730   Hamon PET. A similar signal is also seen for the MC-LSTM-PET but not the MC-LSTM or LSTM,

731   although the LSTM predicts changes in FHV closest to the process models.

732

733   For COM, the process models show a wide range of variability in projected change across sites, from no

734   change to 60 days earlier. For the DL models the range of change is much narrower, and the median change

735   in COM is approximately a week less than the median change across the process models. The earlier shift

736   in COM across all models is consistent with anticipated changes to snow accumulation and melt dynamics

737   under warming, with more water entering the stream during the winter and early spring as precipitation

738   shifts more towards rainfall and snowpack melts off earlier in the year (Byun and Hamlet, 2018; Mote et

739   al., 2018; Kayastha  et al., 2022). However, this effect is seen more dramatically in the process models, as

740   evidenced by more prominent changes to their daily and monthly hydrographs under warming during the

741   winter and early spring as compared to the DL models (see Figures S5 and S6). The method of PET

742   estimation has relatively little impact on both process model and DL based estimates of change in COM.

743

744   We note that the results above do not change even when considering the parametric uncertainty in the

745   process models, although for some metrics (FLV), uncertainty in process model estimated changes due to

746   parametric uncertainty is large (see Figure S7). We also note that if the static watershed properties

747   (pet_mean, aridity, t_mean, frac_snow; see Table 1) are changed to reflect warmer temperatures and higher

748   PET, all three DL models exhibit unrealistic water gains for between 15%-40% of locations depending on

749　the model and PET method, with the most water gains occurring under the LSTM (Figure S8). These results

750　suggest that changing the static watershed properties associated with long-term climate characteristics can

751　degrade the quality of the estimated responses, at least when the temperature shifts are large and the range

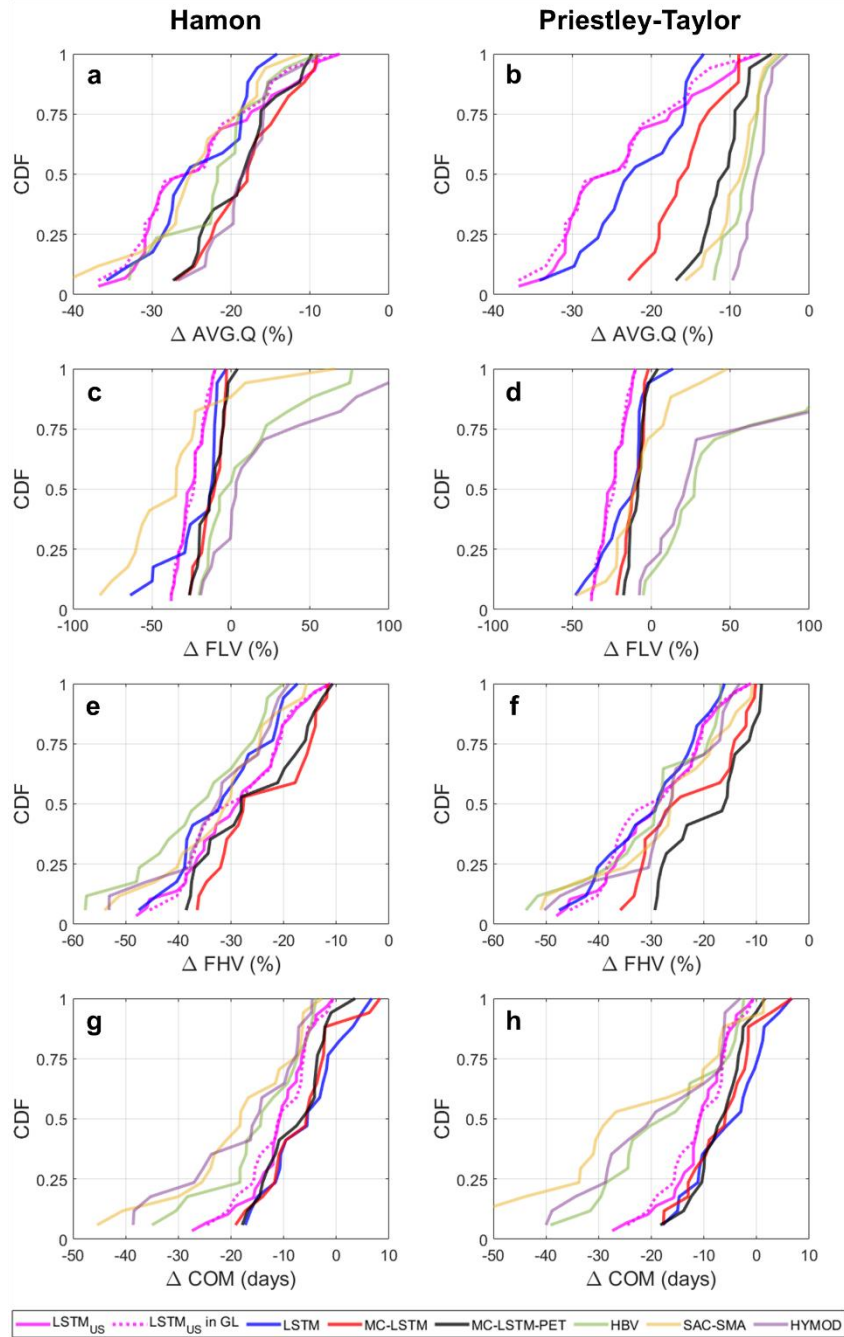752　of average temperature and PET in the training set is limited.
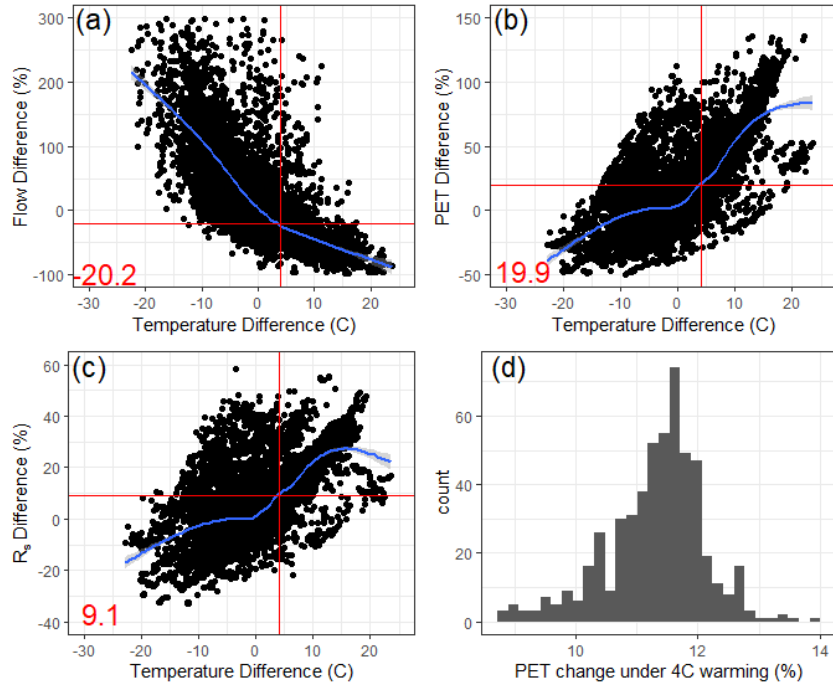
753



754

**Figure 7.** The distribution of change in (a,b) long term mean daily flow (AVG.Q), (c,d) low flows (FLV), (e,f) high flows (FHV), and (g,h) seasonal streamflow timing (COM) across the 141 training sites and all models under a scenario of 4°C warming using (a,c,e,g) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the deep learning models, changes were only made to the dynamic inputs (i.e., no changes to static inputs).

One reason why the Great Lakes LSTM exhibits excessive water losses under warming could be that the model was trained using sites that are confined to a limited range of temperature and PET values found in the Great Lakes basin (spanning approximately 40.5°-50°N), and so is ill-suited to extrapolate hydrologic response under warming conditions that extend beyond this temperature and PET range. To evaluate this hypothesis, we examine changes to AVG.Q, FLV, FHV, and COM under 4°C warming at the 29 CAMELS watersheds within the Great Lakes basin using the National LSTM (Figure 8). For comparison, we also examine similar changes under all six Great Lakes DL and process models at 17 of those 29 CAMELS basins that were used in the training and testing sets for the Great Lakes models. We also highlight the National LSTM predictions for those 17 sites. Note that in Figure 8, the National LSTM predictions do not differ between Hamon and Priestley Taylor PET, because PET is not an input to that model.

The National LSTM was trained to watersheds across the CONUS (spanning approximately 26°-49°N), and so was exposed to watersheds with much warmer conditions and higher PET during training. However, we find that the National LSTM still predicts very large declines in AVG.Q. For the 29 CAMELS watersheds in the Great Lakes basin, the median decline in AVG.Q under the National LSTM is approximately 25%, which is only 0-6% larger than the median predictions of loss under the process models using Hamon PET but 16-19% larger than the process model losses under Priestley-Taylor PET (Figure 8a,b). We also see larger declines in FLV under the National LSTM as compared to the other Great Lakes DL models (Figure 8c,d). The National LSTM predicts changes in FHV (Figure 8e,f) and COM (Figure 8g,h) that are relatively similar to the process models. For COM, the predictions of change are still smaller than the process models but closer to the process models than any Great Lakes DL model, suggesting that the National LSTM predicts shifting snow accumulation and melt dynamics more consistently with the

36

783    process models than regionally fit DL models. In addition, the hydrologic predictions are stable under the

784    National LSTM regardless of whether only dynamic inputs or both dynamic and static inputs are changed

785    under warming (see Figure S9), in contrast to the Great Lakes DL models. Therefore, the use of more

786    watersheds in training than span a more diverse set of climate conditions likely benefit the model when

787    inputs are shifted to reflect new climate conditions. However, as shown in Figure 8a,b, this benefit does not

788    mitigate the tendency for the National LSTM to overestimate water loss under warming.

789

**Figure 8.** The distribution of change in (a,b) long term mean daily flow (AVG.Q), (c,d) low flows (FLV), (e,f) high flows (FHV), and (g,h) seasonal streamflow timing (COM) across 29 CAMELS sites within the Great Lakes basin under the National LSTM (solid pink), as well as for 17 of those 29 sites from the Great Lakes deep learning and process models, under a scenario of 4°C warming. Results from the National LSTM for those 17 sites are also highlighted (dashed pink). For the Great Lakes models only, results differ when using (a,c,e,f) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the National LSTM, changes were made only to the dynamic inputs.

799     To better understand why the National LSTM predicts large water losses under warming, it is instructive

800     to examine how long-term mean streamflow, (Priestly-Taylor estimated) PET, and $R_s$ vary across all 531

801     CAMELS watersheds of different average temperatures, and compare this variability to predicted changes

802     in PET at each site under warming. Specifically, we calculate the difference in long-term (1980-2014) mean

803     streamflow (Figure 9a), PET (Figure 9b), and $R_s$ (Figure 9c) across all pairs of basins in the CAMELS

804     dataset with average long-term precipitation within 1% of each other (i.e., we only examine pairs of basins

805     with very similar long-term mean precipitation). Then, for each basin pair, we plot the difference in long-

806     term mean streamflow, PET, and $R_s$ against the difference in long-term average temperature for that pair.

807     The results show that the difference in long-term mean streamflow across watersheds with similar

808     precipitation becomes negative when the difference in temperature is positive (i.e., warmer watersheds have

809     less flow on average), and that when the difference in average temperature reaches 4°C, flows differ by

810     about 20% on average (Figure 9a). This is very similar to the predicted median decline in long-term mean

811     streamflow seen for the National LSTM in Figure 8. We also note that average PET increases by

812     approximately 20% between watersheds that differ in average temperature by 4°C (Figure 9b). However,

813     higher PET in warmer watersheds is related both to the direct effect of temperature on vapor pressure deficit,

814     as well as to the fact that higher incoming solar radiation co-occurs in warmer watersheds ($R_s$ is

815     approximately 9% higher across watershed pairs that differ by 4°C; Figure 9c). Using the Priestley-Taylor

816     method, we estimate that average PET would only increase by between 9-14% (median of 11.5%) if

817     temperatures warm by 4°C and $R_s$ is held at historic values, while $R_n$ is increased slightly due to declines

818     in net outgoing longwave radiation with warming (Figure 9d). However, the National LSTM appears to

819     convolute the effects of temperature and $R_s$ and cannot separate out their effects on evaporative water loss,

820     leading to larger predicted streamflow losses under 4°C warming than changes in PET would warrant. This

821     is possibly because of the very strong correlation between at-site daily temperature and $R_s$ historically

822     (median correlation of 0.85 across all CAMELS watersheds).

823

**Figure 9.** The percent difference in long-term (1980-2014) mean (a) streamflow, (b) Priestley-Taylor based PET, and (c) downward shortwave radiation ($R_s$) for all pairs of CAMELS basins with average precipitation within 1% of each other, plotted against differences in average temperature for each pair. A loess smooth is provided for each scatter (blue), along with the changes in variable estimated at a 4°C temperature difference between pairs of sites (red). (d) The projected change in Priestley-Taylor based PET (as a percentage) for each CAMELS basin under 4°C warming, assuming no change in $R_s$.

## 5. Discussion and Conclusion

In this study, we contribute a sensitivity analysis that evaluates the physical plausibility of streamflow responses under warming using DL rainfall-runoff models. The basis for this evaluation is anchored to the assumption that differences in estimated streamflow responses should emerge under very different scenarios of PET under warming, and that realistic predictions of PET and water loss under warming tend to be much lower than those estimated by temperature-based PET methods. Accordingly, we assume that physically plausible streamflow predictions should be able to respond to lower energy-budget based PET projections under warming and, all else equal, estimate smaller streamflow losses.

The results of this study show that a standard LSTM is not able to predict physically realistic differences in streamflow response across substantially different estimates of PET under warming. This discrepancy

40

843     emerged despite the fact that the standard LSTM was a far better model for streamflow estimation in

844     ungauged basins compared to three process-based models under historic climate conditions. In addition,

845     the National LSTM trained to a much larger set of watersheds (531 basins across 23° of latitude) using

846     temperature, vapor pressure, and $R_s$ directly (rather than PET) also estimated water loss under warming that

847     far exceeded the losses estimated with process models forced with energy budget-based PET. Since water

848     losses estimated using energy budget-based PET are generally considered more realistic (Lofgren et al.,

849     2011; Shaw and Riha, 2011; Lofgren and Rouhana, 2016; Milly and Dunne, 2017; Lemaitre-Basset et al.

850     2022), this result casts doubt over the physical plausibility of the LSTM predictions.

851

852     Results from this work also suggest that PIML-based DL models can capture physically plausible

853     streamflow responses under warming while still maintaining superior prediction skill compared to process

854     models, at least in some cases. In particular, a mass conserving LSTM that also respected the limits of water

855     loss due to evapotranspiration (the MC-LSTM-PET) was able to predict changes in long-term mean

856     streamflow that much more closely aligned with process-model based estimates, while also providing

857     competitive out-of-sample performance across all models considered (including the other DL models). A

858     more conventional MC-LSTM that did not limit water losses by PET was less consistent with process-based

859     estimates of change in long-term mean streamflow. These results highlight the potential for PIML-based

860     DL models to help achieve similar performance improvements over process-based models as documented

861     in recent work on DL rainfall-runoff models (Kratzert et al., 2019a,b; Feng et al., 2020; Nearing et al., 2021)

862     while also producing projections under climate change that are more consistent with theory than non-PIML

863     DL models.

864

865     An interesting result from this study was the disagreement in the change in high flows and seasonal

866     streamflow timing between all Great Lakes DL models and process models, the latter which estimated

867     greater reductions in high flows and larger shifts of water towards earlier in the year. Predictions from the

868     Great Lakes DL models were also unstable if static climate properties of each watershed were changed

869     under warming. In contrast, the National LSTM was more stable if static properties were changed, and it

870     predicted changes to high flows and seasonal timing that were more like the process models than predictions

871     from the Great Lakes DL models. The results for COM in particular suggest that the National LSTM may

872     be more consistent with the process models in terms of its representation of warming effects on snow

873     accumulation and melt processes and the resulting shifts in the seasonal hydrograph, although differences

874     with the process model predictions were still notable. Still, these results are consistent with past work

875     showing that large-sample LSTMs can learn to represent snow processes internally from meteorological

876     and streamflow data (Lees et al., 2022). While its challenging to know which set of predictions are correct

877     for these streamflow properties, these results overall favor predictions from the National LSTM over the

878     regional LSTMs and highlight the benefits of DL rainfall-runoff models trained to a larger set of diverse

879     watersheds for climate change analysis.

880

881  To properly interpret the results of this work, there are several limitations of this study that require

882  discussion. First there were differences in the inputs and data sources between the National LSTM and all

883  other Great Lakes models, including the source of meteorological data and the lack of PET as an input into

884  the National LSTM.  While this latter discrepancy might be less impactful (i.e., the National LSTM was

885  provided meteorological inputs that together completely determine Hamon and Priestley-Taylor PET), the

886  difference in meteorological data across the two sets of models is a substantial source of uncertainty and

887  could lead to non-trivial differences in hydrologic response estimation, complicating a direct comparison

888  of the National LSTM to the other models. Future work for the Great Lakes Intercomparison Project should

889  consider developing consistent datasets with other (and larger) benchmark datasets like CAMELS to

890  address this issue.

891

892  Another important limitation is how we constructed the warming scenarios, with 4°C warming and shifts

893  to PET but no changes to other meteorological variables (net incoming shortwave radiation, precipitation,

894  humidity, air pressure, wind speeds). These scenarios and associated sensitivity analyses were constructed

895    in the style of other metamorphic tests for hydrologic models (Yang and Chui, 2021; Razavi, 2021; Reichert

896    et al., 2023), where we define input changes with expected responses and test whether model behavior is

897    consistent with these expectations. However, for DL and other machine learning models, the results of such

898    sensitivity analyses may be unreliable because of distributional shifts between the training and testing data

899    and poor out-of-distribution generalization (see Shen et al., 2021, Wang et al., 2023, and references within).

900    When trained, conventional machine learning models try to leverage all of the correlations within the

901    training set to minimize training errors, which is effective in out-of-sample performance only if those same

902    patterns of correlation persistent into the testing data (Liu et al., 2021). In our experimental design, we

903    impose a distinct shift in the joint distribution of the inputs (i.e., a covariate shift) by increasing temperatures

904    and PET but leaving unchanged other meteorological inputs, thereby altering the correlation among inputs.

905    Therefore, one might expect some degradation in the DL model-based predictions of streamflow under

906    these scenarios.

907

908    The challenge of out-of-distribution generalization and its application to DL rainfall-runoff model testing

909    under climate change highlights several important avenues for future work. First, additional efforts are

910    needed to evaluate the physical plausibility of DL-based hydrologic projections under climate change while

911    ensuring that the joint distribution of all meteorological inputs used in future scenarios is realistic. For

912    example, there are physical relationships between changes in temperature and net radiation (Nordling et al.,

913    2021), as well as temperature, humidity, and extreme precipitation (Ali et al., 2018; Najibi et al., 2022),

914    that should all be preserved in future climate scenarios. The use of climate model output may be well suited

915    for such tests, although care is needed to avoid statistical bias correction and downscaling (i.e., post-

916    processing) of multiple climate fields that could cause shifts in the joint distribution across inputs (Maraun,

917    2016). High-resolution convective-permitting models may be helpful in this regard, given their improved

918    accuracy for key climate fields like precipitation (Kendon et al. 2017).

919

43

920   There are also several emerging techniques in machine learning to address out-of-distribution

921   generalization directly. One set of promising methods is causal learning, defined broadly as methods aimed

922   at identifying input variables that have a causal relationship with the target variable and to leverage those

923   inputs for prediction (Shen et al., 2021). PIML approaches, such as the MC-LSTM-PET model proposed

924   in this work, fall into this category (Vasudevan et al., 2021). Here, prior scientific knowledge on casual

925   structures can be embedded into the DL model through tailored loss functions or, as in the case of the MC-

926   LSTM-PET model, through architectural adjustments or constraints (for other examples outside of

927   hydrology, see Lin et al., 2017; Ma et al., 2018). The MC-LSTM-PET model can be viewed as a specific,

928   limited case of a broader class of learnable, differentiable, process-based models (also referred to as hybrid

929   differentiable models; Jiang et al., 2020; Feng et al., 2022; Feng et al., 2023a). These models use process

930   model architectures as a backbone for model structure, which is then enhanced through flexible, data-driven

931   learning for a subset of processes. Recent work has shown that these models can achieve similar

932   performance to LSTMs but can also represent and output different internal hydrologic fluxes (Feng et al.,

933   2022; Feng et al., 2023a).

934

935   However, challenges can arise when imposing architectural constraints in PIML models. For example, the

936   MC-LSTM-PET model makes the assumption that all water loss in the system is due to evapotranspiration,

937   and therefore cannot exceed PET. However, other terminal sinks are possible, such as human water

938   extractions and inter-basin transfers (Siddik et al. 2023) or water lost to aquifer recharge and inter-basin

939   groundwater fluxes (Safeeq et al., 2021; Jasechko et al., 2021). It is difficult to know the magnitude of these

940   alternative sinks given unknown systematic errors in other inputs (e.g., underestimation of precipitation

941   from under-catch) that confound water balance closure analyses. Still, recent techniques and datasets to

942   help quantify these sinks (Gordon et al., 2022; Siddik et al. 2023) provide an avenue to integrate them into

943   the MC-LSTM-PET constraints. Yet as constraints are added to the model architecture, the potential grows

944   for inductive bias that negatively impacts generalizability. For instance, a recent evaluation of hybrid

945   differentiable models showed that they underperformed relative to a standard LSTM due to structural

946  deficiencies in cold regions, arid regions, and basins with considerable anthropogenic impacts (Feng et al.,

947  2023b). Some of these challenges may be difficult to address because only differentiable process models

948  can be considered in this hybrid framework, limiting the process model structures that could be adapted

949  with this approach. Additional work is needed to evaluate the benefits and drawbacks of these different

950  PIML-based approaches, preferably on large benchmarking datasets such as CAMELS or CAVARAN

951  (Kratzert et al., 2023).

952

953  Given some of the potential challenges above, other DL methods that advance causality while making fewer

954  assumptions on watershed-scale process controls are also worth pursuing. For example, a series of

955  techniques have emerged that embed the concept and constraints of directed acyclic graphs within deep

956  neural networks in such a way that the architecture of the neural network is inferred from the data to encode

957  causality among variables (see Luo et al., 2020 and references within). That is, frameworks to optimize the

958  architecture of the model can be designed not only to maximize out-of-sample predictive performance, but

959  also to promote causality. Alternatively, domain-invariant learning attempts to promote the identification

960  of features that are domain-specific versus domain invariant, by separating and labeling training data from

961  different 'domains' or 'environments' (Ilse et al., 2021). In the case of DL rainfall-runoff models, this

962  strategy could be implemented, for instance, by pairing observed climate and streamflow (one domain) with

963  land surface model-based streamflow estimated using future projected climate model output (another

964  domain), with the goal to learn invariant relationships between key climate inputs (e.g., net radiation or

965  PET) and streamflow across the two domains. Here, there may be a benefit from including data from the

966  land surface and climate models, where the correlation between temperature, net radiation, and PET may

967  be weaker under projected climate change. These techniques offer an intriguing alternative for the next

968  generation of DL hydrologic models that can generalize well under climate change, and should be the focus

969  of further exploration.

970

974  **Data Availability Statement**

975  The code used for this project is available at https://doi.org/10.5281/zenodo.10027355. All data used to

976  train and evaluate the models are available at https://doi.org/10.20383/103.0598.

977

978  **References**

979  Ali, H., Fowler, H. J., & Mishra, V. (2018). Global observational evidence of strong linkage between dew
980  point temperature and precipitation extremes. Geophysical Research Letters, 45, 12320–
981  12330. https://doi.org/10.1029/2018gl080557
982
983  Allen, R.G., Pereira, L.S., Raes, D., et al. (1998) Crop Evapotranspiration-Guidelines for Computing
984  Crop Water Requirements-FAO Irrigation and Drainage Paper 56. FAO, Rome, 300(9): D05109.
985
986  Anderson, E. A. (1976). A point energy and mass balance model of a snow cover (NOAA Technical
987  Report NWS 19). Silver Spring, MD: National Oceanic and Atmosphere Administration.
988
989  Bastola S., Murphy C., Sweeney J. (2011). The role of hydrological modelling uncertainties in climate
990  change impact assessments of Irish river catchments. Adv Water Resour., 34, 562–76.
991
992  Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Miralles, D. G., McVicar, T. R., Schellekens, J.,
993  and Bruijnzeel, L. A. (2016), Global-scale regionalization of hydrologic model parameters, Water Resour.
994  Res., 52, 3599–3622, doi:10.1002/2015WR018247.

995  Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Dutra, E., Fink, G., Orth, R., and Schellekens, J.: Global
996  evaluation of runoff from 10 state-of-the-art hydrological models (2017), Hydrol. Earth Syst. Sci., 21,
997  2881–2903, https://doi.org/10.5194/hess-21-2881-2017.

998  Bergström, S. & Forsman, A. (1973) Development of a conceptual deterministic rainfall-runoff model.
999  Nordic Hydrol. 4, 147–170.
1000
1001  Beven, K. (2023). Benchmarking hydrological models for an uncertain future. Hydrological
1002  Processes, 37( 5), e14882. https://doi.org/10.1002/hyp.14882
1003
1004  Boyle, D. P. (2001). Multicriteria calibration of hydrologic models, (Doctoral dissertation). Retrieved from
1005  UA Campus Repository (http://hdl.handle.net/10150/290657), Tucson, AZ: The University of Arizona.
1006
1007  Breuer, L., Huisman, J. A., Willems, P., Bormann, H., Bronstert, A., Croke, B. F. W., Frede, H.-G., Gräff,
1008  T., Hubrechts, L., Jakeman, A. J., Kite, G., Lanini, J., Leavesley, G., Lettenmaier, D. P., Lindström, G.,
1009  Seibert, J., Sivapalan, M., and Viney, N. R.: Assessing the impact of land use change on hydrology by

1010 ensemble modeling (LUCHEM). I: Model intercomparison with current land use, Adv. Water Resour.,
1011 32, 129–146, https://doi.org/10.1016/j.advwatres.2008.10.003, 2009.
1012
1013 Burnash, R. J. (1995). The NWS river forecast system - catchment modeling. In Singh, V. (Ed.), Computer
1014 Models of Watershed Hydrology (pp. 311-366). Littleton, CO: Water Resources Publication.
1015
1016 Byun, K. and Hamlet, A.F. (2018), Projected changes in future climate over the Midwest and Great Lakes
1017 region using downscaled CMIP5 ensembles. Int. J. Climatol, 38: e531-
1018 e553. https://doi.org/10.1002/joc.5388
1019
1020 Campbell, M., Cooper, M. J. P., Friedman, K., & Anderson, W. P. (2015). The economy as a driver of
1021 change in the Great Lakes - St. Lawrence basin. *Journal of Great Lakes Research*, *41*, 69–83.
1022
1023 Clark, M. P., Bierkens, M. F. P., Samaniego, L., Woods, R. A., Uijlenhoet, R., Bennett, K. E., Pauwels,
1024 V. R. N., Cai, X., Wood, A. W., and Peters-Lidard, C. D. (2017). The evolution of process-based
1025 hydrologic models: historical challenges and the collective quest for physical realism, Hydrol. Earth Syst.
1026 Sci., 21, 3427–3440, https://doi.org/10.5194/hess-21-3427-2017.

1027 Clark, M.P., Wilby, R.L., Gutmann, E.D. et al. Characterizing Uncertainty of the Hydrologic Impacts of
1028 Climate Change. Curr Clim Change Rep 2, 55–64 (2016). https://doi.org/10.1007/s40641-016-0034-x
1029
1030 Demargne, J. et al. (2014). The Science of NOAA's Operational Hydrologic Ensemble Forecast
1031 Service. Bull. Amer. Meteor. Soc., 95, 79–98, https://doi.org/10.1175/BAMS-D-12-00081.1.
1032
1033 Fan, Y. (2019). Are catchments leaky? WIREs Water, 6(6). https://doi.org/10.1002/wat2.1386
1034
1035 Feng, D., Fang, K., & Shen, C. (2020). Enhancing streamflow forecast and extracting insights using long-
1036 short term memory networks with data integration at continental scales. Water Resources Research, 56,
1037 e2019WR026793. https://doi.org/ 10.1029/2019WR026793
1038
1039 Feng, D., Liu, J., Lawson, K., & Shen, C. (2022). Differentiable, learnable, regionalized process-based
1040 models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy. Water
1041 Resources Research, 58, e2022WR032404. https://doi.org/10.1029/2022WR032404
1042
1043 Feng, D., Beck, H., Lawson, K., and Shen, C. (2023a). The suitability of differentiable, physics-informed
1044 machine learning hydrologic models for ungauged regions and climate change impact assessment,
1045 Hydrol. Earth Syst. Sci., 27, 2357–2373, https://doi.org/10.5194/hess-27-2357-2023.
1046
1047 Feng, D., Beck, H., de Bruijn, J., Sahu, R. K., Satoh, Y., Wada, Y., Liu, J., Pan, M., Lawson, K., and
1048 Shen, C. (2023b). Deep Dive into Global Hydrologic Simulations: Harnessing the Power of Deep
1049 Learning and Physics-informed Differentiable Models (δHBV-globe1.0-hydroDL), Geosci. Model Dev.
1050 Discuss. [preprint], https://doi.org/10.5194/gmd-2023-190, in review.
1051
1052 Frame, J.M., Kratzert, F., Gupta, H.V., Ullrich, P., & Nearing, G.S. (2022). On Strictly enforced mass
1053 conservation constraints for modeling the Rainfall-Runoff process. Hydrological Processes, 37, e14847,
1054 https://doi.org/10.1002/hyp.14847.
1055
1056 Frame, J.M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., et al. (2021b). Deep learning
1057 rainfall-runoff predictions of extreme events. *Hydrology and Earth System Sciences*, 26, 3377-
1058 3392, https://doi.org/10.5194/hess-26-3377-2022.
1059

Frame, J.M., Kratzert, F., Raney II, A., Rahman, M., Salas, F.R., & Nearing, G.S. (2021a). Post-processing the National Water Model with Long Short-Term Memory networks for streamflow predictions and diagnostics. *Journal of the American Water Resources Association*, 1-12. https://doi.org/10.1111/1752-1688.12964

Fry, L. M., Hunter, T. S., Phanikumar, M. S., Fortin, V., and Gronewold, A. D. (2013), Identifying streamgage networks for maximizing the effectiveness of regional water balance modeling, Water Resour. Res., 49, 2689– 2700, doi:10.1002/wrcr.20233.

Gasset, N., Fortin, V., Dimitrijevic, M., Carrera, M., Bilodeau, B., Muncaster, R., Gaborit, É., Roy, G., Pentcheva, N., Bulat, M., Wang, X., Pavlovic, R., Lespinas, F., Khedhaouiria, D., and Mai, J.: A 10 km North American precipitation and land-surface reanalysis based on the GEM atmospheric model, Hydrol. Earth Syst. Sci., 25, 4917–4945, https://doi.org/10.5194/hess-25-4917-2021, 2021.

Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., & Hochreiter, S. (2021a). Rainfall-runoff prediction at multiple timescales with a single Long Short-Term Memory network. *Hydrology and Earth System Sciences*, *25*, 2045-2062. https://doi.org/10.5194/hess-25-2045-2021

Gauch, M., Mai, J., & Lin, J. (2021b). The proper care and feeding of CAMELS: How limited training data affects streamflow prediction. *Environmental Modelling and Software*, *135*, 104926. https://doi.org/10.1016/j.envsoft.2020.104926

Gordon, B. L., Crow, W. T., Konings, A. G., Dralle, D. N., & Harpold, A. A. (2022). Can we use the water budget to infer upland catchment behavior? The role of data set error estimation and interbasin groundwater flow. *Water Resources Research*, *58*, e2021WR030966. https://doi.org/10.1029/2021WR030966

Greve, P., Roderick, M.L., Ukkola, A.M., and Wada, Y. (2019), The aridity index under global warming, Environmental Research Letters, 14, 124006, https://doi.org/10.1088/1748-9326/ab5046.

Gronewold, A. D., and Rood, R. B. (2019). Recent water level changes across Earth's largest lake system and implications for future variability. *Journal of Great Lakes Research*, *45*(1), 1–3.

Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F. (2009). Decom- position of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling, J. Hydrol., 377, 80–91.

Hamon, W. R. (1963). Estimating Potential Evapotranspiration, T. Am. Soc. Civ. Eng., 128, 324–338, https://doi.org/10.1061/TACEAT.0008673.

Hansen, C., Shafiei Shiva, J., McDonald, S., and Nabors, A. (2019). Assessing Retrospective National Water Model Streamflow with Respect to Droughts and Low Flows in the Colorado River Basin. Journal of the American Water Resources Association 964– 975. https://doi.org/10.1111/1752-1688.12784.

Hargreaves, G.H., and Samani, Z.A. (1985). Reference crop evapotranspiration from temperature. Applied Engineering in Agriculture 1: 96–99.

Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735-1780. https://doi.org/10.1162/neco.1997.9.8.1735

Hoedt, P.J., F. Kratzert, D. Klotz, C. Halmich, M. Holzleitner, G. Nearing, et al. (2021). MC-LSTM:

Mass-Conserving LSTM. *arXiv e-prints*, arXiv:2101.05186. Retrieved from
https://arxiv.org/abs/2101.05186


Höge, M., Scheidegger, A., Baity-Jesi, M., Albert, C., and Fenicia, F. (2022). Improving hydrologic
models for predictions and process understanding using neural ODEs, Hydrol. Earth Syst. Sci., 26, 5085–
5102, https://doi.org/10.5194/hess-26-5085-2022.


Hrachowitz, M. et al. (2013). A decade of Predictions in Ungauged Basins (PUB)—a
review, Hydrological Sciences Journal, 58:6, 1198-1255, DOI: 10.1080/02626667.2013.803183


Ilse, M., Tomczak, J.M., and Forré, P. (2021). Selecting Data Augmentation for Simulating Interventions.
Proceedings of the 38th International Conference on Machine Learning, PMLR 139:4555-4562.


Jasechko, S., Seybold, H., Perrone, D. et al. Widespread potential loss of streamflow into underlying
aquifers across the USA. Nature 591, 391–395 (2021). https://doi.org/10.1038/s41586-021-03311-x


Jiang, S., Zheng, Y., & Solomatine, D. (2020). Improving AI system awareness of geoscience knowledge:
Symbiotic integration of physical approaches and deep learning. Geophysical Research Letters, 46,
e2020GL088229. https://doi. org/10.1029/2020GL088229


Karpantne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., et al. (2017).
Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Transactions on
Knowledge and Data Engineering*, *29*(10), 2318-2331. https://doi.org/10.1109/TKDE.2017.2720168


Kayastha, M.B., Ye, X., Huang, C., and Xue, P. (2022), Future rise of the Great Lakes water levels under
climate change, Journal of Hydrology, 612 (Part B), 128205,
https://doi.org/10.1016/j.jhydrol.2022.128205.


Kendon, Elizabeth J., Nikolina Ban, Nigel M. Roberts, Hayley J. Fowler, Malcolm J. Roberts, Steven C.
Chan, Jason P. Evans, Giorgia Fosser, and Jonathan M. Wilkinson. (2017). Do Convection-Permitting
Regional Climate Models Improve Projections of Future Precipitation Change? *Bulletin of the American
Meteorological Society* 98 (1): 79–93. https://doi.org/10.1175/BAMS-D-15-0004.1.


Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *arXiv e-prints*,
arXiv:1412.6980. Retrieved from https://arxiv.org/abs/1412.6980


Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S.,
and Nearing, G. (2022). Uncertainty estimation with deep learning for rainfall–runoff modeling, Hydrol.
Earth Syst. Sci., 26, 1673–1693, https://doi.org/10.5194/hess-26-1673-2022.


Konapala, G., Kao, S. C., Painter, S., & Lu, D. (2020). Machine learning assisted hybrid models can
improve streamflow simulation in diverse catchments across the conterminous US. Environmental
Research Letters, 15(10), 104022. https://doi.org/10.1088/1748-9326/aba927


Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019a).
Toward improved predictions in ungauged basins: Exploiting the power of machine learning. *Water
Resources Research*, *55*, 11,344–11,354. https://doi.org/10.1029/2019WR026065


Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. S. (2019b). Towards
learning universal, regional, and local hydrological behaviors via machine learning applied to large-

sample datasets. *Hydrology and Earth System Sciences*, *23*, 5089-5110. https://doi.org/10.5194/hess-23-5089-2019

Kratzert, F., Klotz, D., Hochreiter, S., & Nearing, G. S. (2021). A note on leveraging in multiple meteorological data sets with deep learning for rainfall-runoff modeling. Hydrology and Earth System Sciences, 25(5), 2685–2703. https://doi.org/10.5194/hess-25-2685-2021.

Kratzert, F., Nearing, G., Addor, N. et al. (2023), Caravan - A global community dataset for large-sample hydrology. Sci Data 10, 61. https://doi.org/10.1038/s41597-023-01975-w

Krøgli, I. K., Devoli, G., Colleuille, H., Boje, S., Sund, M., and Engen, I. K.: The Norwegian forecasting and warning service for rainfall- and snowmelt-induced landslides, Nat. Hazards Earth Syst. Sci., 18, 1427–1450, https://doi.org/10.5194/nhess-18-1427-2018, 2018.

Krysanova, V., Donnelly, C., Gelfan, A., Gerten, D., Arheimer, B., Hattermann, F. and Kundzewicz Z.W. (2018) How the performance of hydrological models relates to credibility of projections under climate change, Hydrological Sciences Journal, 63:5, 696-720, DOI: 10.1080/02626667.2018.1446214

Lai, C., Chen, X., Zhong, R., and Wang, Z. (2022), Implication of climate variable selections on the uncertainty of reference crop evapotranspiration projections propagated from climate variables projections under climate change, Agricultural Water Management, 259(1), 107273, https://doi.org/10.1016/j.agwat.2021.107273.

Lee, D., Lee, G., Kim, S., & Jung, S. (2020). Future Runoff Analysis in the Mekong River Basin under a Climate Change Scenario Using Deep Learning. *Water*, *12*(6):1556. https://doi.org/10.3390/w12061556

Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve, P., Slater, L., and Dadson, S. J. (2022). Hydrological concept formation inside long short-term memory (LSTM) networks, Hydrol. Earth Syst. Sci., 26, 3079–3101, https://doi.org/10.5194/hess-26-3079-2022.

Lehner, B., Verdin, K., and Jarvis, A. (2008). New Global Hydrography Derived From Spaceborne Elevation Data, Eos T. Am. Geophys. Un., 89, 93–94.

Lemaitre-Basset, T., Oudin, L., Thirel, G., and Collet, L.: Unraveling the contribution of potential evaporation formulation to uncertainty under climate change, Hydrol. Earth Syst. Sci., 26, 2147–2159, https://doi.org/10.5194/hess-26-2147-2022, 2022.

Li, K., Huang, G., Wang, S., Razavi, S., & Zhang, X. (2022). Development of a joint probabilistic rainfall-runoff model for high-to-extreme flow projections under changing climatic conditions. Water Resources Research, 58, e2021WR031557. https://doi. org/10.1029/2021WR031557

Lin, L., Gettelman, A., Fu, Q. et al. Simulated differences in 21st century aridity due to different scenarios of greenhouse gases and aerosols. Climatic Change 146, 407–422 (2018). https://doi.org/10.1007/s10584-016-1615-3

Lin, C., Jain, S., Kim, H., Bar-Joseph, Z. (2017). Using neural networks for reducing the dimensions of single-cell RNA-Seq data, Nucleic Acids Research, Volume 45, Issue 17, 29 September 2017, Page e156, https://doi.org/10.1093/nar/gkx681

Liu, J., Hu, Z., Cui, P., Li, B., and Shen, Z. (2021). Heterogeneous risk minimization. In ICML, PMLR. PMLR.

1210
1211  Liu, X., Li, C., Zhao, T., and Han, L. (2020) Future changes of global potential evapotranspiration
1212  simulated from CMIP5 to CMIP6 models, Atmospheric and Oceanic Science Letters, 13:6, 568-
1213  575, DOI: 10.1080/16742834.2020.1824983
1214
1215  Liu, Z., Han, J., and Yang, H. (2022), Assessing the ability of potential evaporation models to capture the
1216  sensitivity to temperature, Agricultural and Forest Meteorology, 317, 108886.
1217
1218  Lofgren, B.M., Hunter, T.S., Wilbarger, J. (2011), Effects of using air temperature as a proxy for potential
1219  evapotranspiration in climate change scenarios of Great Lakes basin hydrology, Journal of Great Lakes
1220  Research, 37 (4), 744-752.
1221
1222  Lofgren, B. M., and Rouhana, J. (2016) Physically Plausible Methods for Projecting Changes in Great
1223  Lakes Water Levels under Climate Change Scenarios. J. Hydrometeor., 17, 2209–
1224  2223, https://doi.org/10.1175/JHM-D-15-0220.1.
1225
1226  Lu, D., Konapala, G., Painter, S. L., Kao, S. C., & Gangrade, S. (2021). Streamflow simulation in data-
1227  scarce basins using Bayesian and physics-informed machine learning models. Journal of
1228  Hydrometeorology, 22(6), 1421– 1438. https://doi.org/10.1175/JHM-D-20-0082.1
1229
1230  Lu, J., Sun, G., McNulty, S.G. and Amatya, D.M. (2005), A comparison of six potential
1231  evapotranspiration methods for regional use in the southeastern United States. JAWRA Journal of the
1232  American Water Resources Association, 41: 621-633. https://doi.org/10.1111/j.1752-
1233  1688.2005.tb03759.x
1234
1235  Luo, Y., Peng, J. & Ma, J. (2020). When causal inference meets deep learning. Nat Mach Intell 2, 426–
1236  427. https://doi.org/10.1038/s42256-020-0218-x
1237
1238  Ma, J., Yu, M., Fong, S. et al. (2018). Using deep learning to model the hierarchical structure and
1239  function of a cell. Nat Methods 15, 290–298. https://doi.org/10.1038/nmeth.4627
1240
1241  Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., et al. (2021). Transferring hydrologic
1242  data across continents – leveraging data-rich regions to improve hydrologic prediction in data-sparse
1243  regions. Water Resources Research, 57, e2020WR028600. https://doi. org/10.1029/2020WR028600
1244
1245  Mai et al. (2022). The Great Lakes runoff intercomparison project phase 4: the Great Lakes (GRIP-GL),
1246  Hydrologic and Earth System Sciences, 26 (13), 3537-3573, https://doi.org/10.5194/hess-26-3537-2022.
1247
1248  Martens, B., Miralles, D. G., Lievens, H., van der Schalie, R., de Jeu, R. A. M., Fernández-Prieto, D.,
1249  Beck, H. E., Dorigo, W. A., and Verhoest, N. E. C. (2017). GLEAM v3: satellite-based land evaporation
1250  and root-zone soil moisture, Geosci. Model Dev., 10, 1903– 1925, https://doi.org/10.5194/gmd-10-1903-
1251  2017.
1252
1253  Maraun, D. (2016). Bias Correcting Climate Change Simulations - a Critical Review. Curr Clim Change
1254  Rep 2, 211–220. https://doi.org/10.1007/s40641-016-0050-x
1255
1256  Melsen, L. A., Addor, N., Mizukami, N., Newman, A. J., Torfs, P. J. J. F., Clark, M. P., Uijlenhoet, R., and
1257  Teuling, A. J. (2018). Mapping (dis)agreement in hydrologic projections, Hydrol. Earth Syst. Sci., 22,
1258  1775–1791, https://doi.org/10.5194/hess-22-1775-2018.
1259

1260  Merz, R., Parajka, J., and Blöschl, G. (2011), Time stability of catchment model parameters: Implications
1261  for climate impact analyses, Water Resour. Res., 47, W02531, doi:10.1029/2010WR009505.

1263  Milly, P.C.D. and Dunne, Krista A. (2017). A Hydrologic Drying Bias in Water-Resource Impact
1264  Analyses of Anthropogenic Climate Change. Journal of the American Water Resources
1265  Association (JAWRA) 53( 4): 822– 838. https://doi.org/10.1111/1752-1688.12538

1267  Mote, P. W., Li, S., Lettenmaier, D. P., Xiao, M., & Engel, R. (2018). Dramatic declines in snowpack in
1268  the western US. *npj Climate and Atmospheric Science*, *1:2*. https://doi.org/10.1038/s41612-018-0012-1

1270  NACLMS: NACLMS website, http://www.cec.org/north-american- environmental-atlas/land-cover-2010-
1271  landsat-30m/ (last access: 31 May 2023), 2017.

1273  Najibi, N., Mukhopadhyay, S., & Steinschneider, S. (2022). Precipitation scaling with temperature in the
1274  Northeast US: Variations by weather regime, season, and precipitation intensity. Geophysical Research
1275  Letters, 49, e2021GL097100. https://doi.org/10.1029/2021GL097100

1277  Nash, J. E. and Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I – A
1278  discussion of principles, J. Hydrol., 10, 282–290.

1280  Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., et al. (2021). What
1281  role does hydrological science play in the age of machine learning? *Water Resources Research*, *57*,
1282  e2020WR028091. https://doi.org/10.1029/2020WR028091

1284  Nearing, G. S., Klotz, D., Frame, J. M., Gauch, M., Gilon, O., Kratzert, F., Sampson, A. K., Shalev, G.,
1285  and Nevo, S. (2022). Technical note: Data assimilation and autoregression for using near-real-time
1286  streamflow observations in long short-term memory networks, Hydrol. Earth Syst. Sci., 26, 5493–5513,
1287  https://doi.org/10.5194/hess-26-5493-2022.

1289  Newman, A., Clark, M. P., Sampson, K., Wood, A., Hay, L., Bock, A., et al. (2015). Development of a
1290  large-sample watershed-scale hydrometeorological dataset for the contiguous USA: Data set
1291  characteristics and assessment of regional variability in hydrologic model performance. *Hydrology and
1292  Earth System Sciences*, *19*(1), 209-223. https://doi.org/10.5194/hess-19-209-2015

1294  Nordling, K., Korhonen, H., Raisanen, J., Partanen, A.-I., Samset, B.H., and Merikanto, J. (2021),
1295  Understanding the surface temperature response and its uncertainty to $CO_2$, $CH_4$, black carbon, and
1296  sulfate, Atmos. Chem. Phys., 21, 14941-14958.

1298  Olsson, J., and Lindstrom, G. (2008), Evaluation and calibration of operational hydrological ensemble
1299  forecasts in Sweden Journal of Hydrology, 350 (1–2), 14-24.

1301  Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andreassian, V., Anctil, F., and Loumagne,
1302  C. (2005). Which potential evapotranspiration input for a lumped rainfall-runoff model? Part 2—Towards
1303  a simple and efficient potential evapotranspiration model for rainfall-runoff modeling. Journal of
1304  Hydrology 303: 290–306.

1306  Plesca, I., Timbe,E., Exbrayat, J.F., Windhorst, D., Kraft, P., Crespo, P., Vachéa, K.B., Frede, H.G., and
1307  Breuer,L. (2012). Model intercomparison to explore catchment functioning: Results from a remote
1308  montane tropical rainforest, Ecol. Model., 239, 3–13.

1310 Priestley, C. H. B., and Taylor, R. J. (1972). On the Assessment of Surface Heat Flux and Evaporation
1311 Using Large-Scale Parameters. Mon. Wea. Rev., 100, 81–92, https://doi.org/10.1175/1520-
1312 0493(1972)100<0081:OTAOSH>2.3.CO;2.

1314 Pryor, S.C., Barthelmie, R.J., Bukovsky, M.S. et al. Climate change impacts on wind power
1315 generation. Nat Rev Earth Environ 1, 627–643 (2020). https://doi.org/10.1038/s43017-020-0101-7

1317 Razavi, S. (2021). Deep learning, explained: Fundamentals, explainability, and bridgeability to process-
1318 based modelling, Environmental Modelling and Software,
1319 105159, https://doi.org/10.1016/j.envsoft.2021.105159.

1321 Reichert, P., Ma, K., Höge, M., Fenicia, F., Baity-Jesi, M., Feng, D., and Shen, C.: Metamorphic Testing
1322 of Machine Learning and Conceptual Hydrologic Models, Hydrol. Earth Syst. Sci. Discuss. [preprint],
1323 https://doi.org/10.5194/hess-2023-168, in review, 2023.

1325 Safeeq, M., Bart, R. R., Pelak, N. F., Singh, C. K., Dralle, D. N., Hartsough, P., & Wagenbrenner, J. W.
1326 (2021). How realistic are water-balance closure assumptions? A demonstration from the southern sierra
1327 critical zone observatory and kings river experimental watersheds. Hydrological Processes, 35: e14199.
1328 https://doi.org/10.1002/hyp.14199

1330 Seibert, J. and Bergström, S. (2022). A retrospective on hydrological catchment modelling based on half a
1331 century with the HBV model, Hydrol. Earth Syst. Sci., 26, 1371–1388, https://doi.org/10.5194/hess-26-
1332 1371-2022.

1334 Shangguan, W., Dai, Y., Duan, Q., Liu, B., and Yuan, H. (2014). A global soil data set for earth system
1335 modeling, J. Adv. Model. Earth Sy., 6, 249–263.

1337 Shaw, S.B. and Riha, S.J. (2011), Assessing temperature-based PET equations under a changing climate
1338 in temperate, deciduous forests. Hydrol. Process., 25: 1466-1478. https://doi.org/10.1002/hyp.7913

1340 Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., and Cui, P. (2021). Towards out-of-distribution
1341 generalization: A survey. arXiv preprint arXiv:2108.13624.

1343 Siddik, M.A.B., Dickson, K.E., Rising, J. et al. Interbasin water transfers in the United States and
1344 Canada. Sci Data 10, 27 (2023). https://doi.org/10.1038/s41597-023-01935-4

1346 Steinman, A.D. et al. (2017), Ecosystem services in the Great Lakes, Journal of Great Lakes Research, 43
1347 (3), 161-168. https://doi.org/10.1016/j.jglr.2017.02.004

1349 Su, Q., & Singh, V. P. (2023). Calibration-free Priestley-Taylor method for reference evapotranspiration
1350 estimation. Water Resources Research, 59, e2022WR033198. https://doi.org/10.1029/2022WR033198

1352 Szilagyi, J., Crago, R., and Qualls, R. (2017), A calibration-free formulation of the complementary
1353 relationship of evaporation for continental-scale hydrology, J. Geophys. Res. Atmos., 122, 264– 278,
1354 doi:10.1002/2016JD025611.

1356 Towler, E., Foks, S. S., Dugger, A. L., Dickinson, J. E., Essaid, H. I., Gochis, D., Viger, R. J., and Zhang,
1357 Y. (2023): Benchmarking high-resolution hydrologic model performance of long-term retrospective
1358 streamflow simulations in the contiguous United States, Hydrol. Earth Syst. Sci., 27, 1809–1825,
1359 https://doi.org/10.5194/hess-27-1809-2023.

1360     Vasudevan, R.K., Ziatdinov, M., Vlcek, L. et al. (2021). Off-the-shelf deep learning is not enough, and
1361     requires parsimony, Bayesianity, and causality. npj Comput Mater 7, 16. https://doi.org/10.1038/s41524-
1362     020-00487-0

1364     Wallner, M., and Haberlandt, U. (2015), Non-stationary hydrological model parameters: a framework
1365     based on SOM-B. Hydrol. Process., 29, 3145–3161. doi: 10.1002/hyp.10430.

1367     Wang, Q. J. (1991). The genetic algorithm and its application to calibrating conceptual rainfall-runoff
1368     models, *Water Resources Research*, *27*(9), 2467-2471. https://doi.org/10.1029/91WR01305

1370     Wang, J., Lan, C., Liu, C., Ouyang, Y., Qin, T., Lu, W., Chen, Y., Zeng, W., Yu, P.S.
1371     (2023). Generalizing to Unseen Domains: A Survey on Domain Generalization, in IEEE Transactions on
1372     Knowledge and Data Engineering, vol. 35, no. 8, pp. 8052-8072, 1 Aug. 2023, doi:
1373     10.1109/TKDE.2022.3178128.

1375     Wi, S., & Steinschneider, S. (2022). Assessing the physical realism of deep learning hydrologic model
1376     projections under climate change. Water Resources Research, 58,
1377     e2022WR032123. https://doi.org/10.1029/2022WR032123

1379     Wu, H., Zhu, W., and Huang, B. (2021), Seasonal variation of evapotranspiration, Priestley-Taylor
1380     coefficient and crop coefficient in diverse landscapes, Geography and Sustainability, 2(3), 224-233,
1381     https://doi.org/10.1016/j.geosus.2021.09.002

1383     Yan, H., Sun, N., Eldardiry, H., Thurber, T. B., Reed, P. M., Malek, K., et al. (2023). Large ensemble
1384     diagnostic evaluation of hydrologic parameter uncertainty in the Community Land Model Version 5
1385     (CLM5). Journal of Advances in Modeling Earth Systems, 15,
1386     e2022MS003312. https://doi.org/10.1029/2022MS003312

1388     Yang, Y., & Chui, T. F. M. (2021). Reliability assessment of machine learning models in hydrological
1389     predictions through metamorphic testing. Water Resources Research, 57,
1390     e2020WR029471. https://doi.org/10.1029/2020WR029471

1392     Yilmaz, K. K., Gupta, H. V., and Wagener, T. (2008). A process-based diagnostic approach to model
1393     evaluation: Application to the NWS distributed hydrologic model, Water Resour. Res., 44, 1–18.

1395     Zhong, L., Lei, H., & Gao, B. (2023). Developing a physics-informed deep learning model to simulate
1396     runoff response to climate change in Alpine catchments. Water Resources Research, 59,
1397     e2022WR034118. https://doi.org/10.1029/2022WR034118