1  # On the need for physical constraints in deep learning rainfall-runoff

2  # projections under climate change

3

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

5  [1]Department of Biological and Environmental Engineering, Cornell University, Ithaca, NY, USA

6  *Correspondence to*: Sungwook Wi (sw2275@cornell.edu)

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26 **Abstract**

27 Deep learning rainfall-runoff models have recently emerged as state-of-the-science tools for hydrologic

28 prediction that outperform conventional, process-based models in a range of applications. However, it

29 remains unclear whether deep learning models can produce physically plausible projections of streamflow

30 under significant amounts of climate change. We investigate this question here, focusing specifically on

31 modeled responses to increases in temperature and potential evapotranspiration (PET). Previous research

32 has shown that temperature-based methods to estimate PET lead to overestimates of water loss in rainfall-

33 runoff models under warming, as compared to energy budget-based PET methods. Consequently, we assess

34 the reliability of streamflow projections under warming by comparing projections with both temperature-

35 based and energy budget-based PET estimates, assuming that reliable streamflow projections should exhibit

36 less water loss when forced with smaller (energy budget-based) projections of future PET. We conduct this

37 assessment using three process-based rainfall-runoff models and three deep learning models, trained and

38 tested across 212 watersheds in the Great Lakes basin. The deep learning models include a regional Long

39 Short-Term Memory network (LSTM), a mass-conserving LSTM (MC-LSTM) that preserves the water

40 balance, and a novel variant of the MC-LSTM that also respects the relationship between PET and water

41 loss (MC-LSTM-PET). We first compare historical streamflow predictions from all models under spatial

42 and temporal validation, and also assess model skill in estimating watershed-scale evapotranspiration. We

43 then force all models with scenarios of warming, historical precipitation, and both temperature-based

44 (Hamon) and energy budget-based (Priestley-Taylor) PET, and compare their projections for changes in

45 average flow, as well as low flows, high flows, and streamflow timing. Finally, we also explore similar

46 projections using a National LSTM fit to a broader set of 531 watersheds across the contiguous United

47 States. The main results of this study are as follows:

48     1. The three Great Lakes deep learning models significantly outperform all process models in

49        streamflow estimation under spatiotemporal validation, with only small differences between the

50      DL models. The MC-LSTM-PET also matches the best process models and outperforms the MC-

51      LSTM in estimating evapotranspiration under spatiotemporal validation.

52   2.  All process models show a downward shift in average flows under warming, but this shift is

53      significantly larger under temperature-based PET estimates than energy budget-based PET. The

54      MC-LSTM-PET model exhibits similar differences in water loss across the different PET forcings,

55      consistent with the process models. However, the LSTM exhibits unrealistically large water losses

56      under warming as compared to the process models using Priestley-Taylor PET, while the MC-

57      LSTM is relatively insensitive to PET method.

58   3.  All deep learning models exhibit smaller changes in high flows and streamflow timing as compared

59      to the process models, while deep learning projections of low flows are all very consistent and

60      within the range projected by process models.

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

62      warming. However, when compared to the Great Lakes deep learning models, projections from the

63      National LSTM were more stable when many inputs were changed under warming and better

64      aligned with process model projections for streamflow timing. This suggests that the addition of

65      more, diverse watersheds in training does help improve climate change projections from deep

66      learning models, but this strategy alone may not guarantee reliable projections under unprecedented

67      climate change.

68   Ultimately, the results of this work suggest that physical considerations regarding model architecture and

69   input variables are necessary to promote the physical realism of deep learning-based hydrologic projections

70   under climate change.

71

72   **Keywords**

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

74   change, rainfall-runoff

## 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 many to argue that DL models represent the state-of-the-science in rainfall-runoff modeling.

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; 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. 2021), suggesting

100    that DL hydrologic models can learn fundamental hydrologic processes. A corollary to this finding is that

101    these models may produce physically plausible streamflow predictions under new climate conditions.

102

103    It is challenging to assess the physical plausibility of DL-based hydrologic projections under significantly

104    different climate conditions, because there are no future observations against which to compare. Recently,

105    Wi and Steinschneider (2022) (hereafter WS22) addressed this challenge directly, forwarding an

106    experimental design in which DL hydrologic models fit to 15 watersheds in California and 531 catchments

107    across the United States were forced with historical precipitation and temperature, but with temperatures

108    adjusted by up to 4°C. Based on past literature (Cayan et al., 2001; Stewart et al., 2005; Kapnick and Hall,

109    2010; Lehner et al., 2017; McCabe et al., 2017; Dierauer et al., 2018; Mote et al., 2018; Woodhouse &

110    Pederson, 2018; Martin et al., 2020; Milly & Dunne, 2020; Rungee et al., 2021; Gordon et al., 2022; Liu et

111    al., 2022), WS22 posited that physically plausible hydrologic projections should show a decline in total

112    annual average streamflow compared to a baseline historical simulation, due to increases in potential

113    evapotranspiration (PET) with warming (and no changes in precipitation). Results showed that the LSTM

114    trained to the 15 watersheds in California often led to misleading increases in annual runoff under

115    significant warming, while this phenomenon was less likely (though still present) in the model trained to

116    531 catchments.

117

118    WS22 also conducted their experiment with physics-informed machine learning (PIML) models, in which

119    data-driven techniques are imbued with process-knowledge constructs (Karpatne et al., 2017). WS22

120    focused on two PIML strategies for the smaller case study in California, using process model output (e.g.,

121    soil moisture, evapotranspiration (ET)) directly as input to the LSTM (similar to Konapala et al., 2020; Lu

122    et al., 2021; Frame et al., 2021a), and also as additional target variables in a multi-output architecture. The

123    former approach had some success in removing instances of increasing runoff ratio with warming, but this

124    depended heavily on the accuracy of the process-model ET.

125

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

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

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

129  LSTM) that ensures cumulative streamflow predictions do not exceed precipitation inputs. This architecture

130  slightly underperformed a standard LSTM when predicting out-of-sample extreme events (Frame et al.,

131  2021b), and some have argued that these physical constraints may inhibit the ability of DL models to learn

132  biases in forcing data (Frame et al. 2022). Still, the benefits of this mass conserving architecture have not

133  been tested when employed under previously unobserved climate change.

134

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

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

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

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

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

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

141  2016; Milly and Dunne, 2017; Lemaitre-Basset et al. 2022). Broadly, this work has shown that temperature-

142  based PET estimation methods (e.g., Hamon, Thornthwaite) significantly overestimate increases in PET

143  under warming as compared to energy budget-based PET estimation methods (e.g., Penman-Monteith,

144  Priestley-Taylor), and consequently lead to unrealistic declines in streamflow under climate change. This

145  is because the actual drying power of the atmosphere is driven by the availability of energy at the surface

146  from net radiation, the current moisture content of the air, temperature (and its effect on the water holding

147  capacity of the air and vapor pressure deficit), and wind speeds. Energy budget-based methods account for

148  some or all of these factors in ways that are generally consistent with their causal impact on PET, while

149  temperature-based methods estimate PET using empirical relationships based largely or entirely on

150  temperature. The latter approach works sufficiently well for rainfall-runoff modeling under historical

151  conditions because of the strong correlation between temperature, net radiation, and PET on seasonal

152    timescales, even though this correlation weakens considerably at shorter timescales (Lofgren et al., 2011).

153    Under climate change, consistent and prominent increases are projected for temperature, but projected

154    changes are less prominent or more uncertain for other factors affecting PET (Lin et al., 2018; Pryor et al.,

155    2020, Liu et al. 2020). Consequently, temperature-based PET methods significantly overestimate future

156    projections of PET compared to energy budget-based methods.

157

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

159    estimation methods under warming provides a unique opportunity to assess the physical plausibility of

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

161    models and forward an experimental design in which both process-based and DL hydrologic models are

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

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

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

165    calculated from the warmed temperatures using both PET estimation methods. We anticipate that the

166    process models 1) will exhibit similar performance in historical training and testing periods when using

167    either temperature-based or energy budget-based PET estimates; but 2) will exhibit significantly larger

168    streamflow declines under warming when using future PET estimated with a temperature-based method. If

169    the DL rainfall-runoff models follow the same pattern, this would suggest that these models are able to

170    learn the role of PET on evaporative water loss. However, if DL-based models estimate similar and large

171    streamflow declines regardless of the method used to estimate and project PET, this would suggest that the

172    DL models did not learn a mapping between PET and water loss. Rather, the DL models learned the

173    historical (but non-causal) correlation between temperature and evaporative water loss, and then incorrectly

174    extrapolated that effect into the future with warmer temperatures. Such an outcome would indicate that

175    some degree of PIML is necessary to guide a DL model towards physically plausible projections under

176    climate change, in contrast to previous arguments against the need for such physical constraints (Frame et

177    al. 2022).

178

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

180    both standard and PIML-based LSTMs. We hypothesize that a standard LSTM will produce unrealistic

181    hydrologic projections because it relies on historical and geographically pervasive correlations between

182    temperature and PET to project streamflow losses under warming. We also hypothesize that PIML-based

183    DL models will be better able to relate future projections of temperature and PET to streamflow change,

184    especially those PIML approaches that directly map PET to evaporative water loss in their architecture.

185

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

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

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

189    2023). The Great Lakes provides an important case study for this work, given their importance to the culture,

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

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

192    in part because of uncertainty in future runoff draining into the lakes from a large contributing area

193    (Kayastha et al. 2022), much of which is ungauged (Fry et al., 2013). Improved rainfall-runoff models that

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

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

196    Intercomparison Project Phase 4 (GRIP-GL) presented in Mai et al. (2022). This study currently provides

197    one of the most robust benchmarks comparing DL rainfall-runoff models to a range of process-based

198    models, and so we design our experiment to be consistent with the data and model development rules

199    outlined in the GRIP-GL.
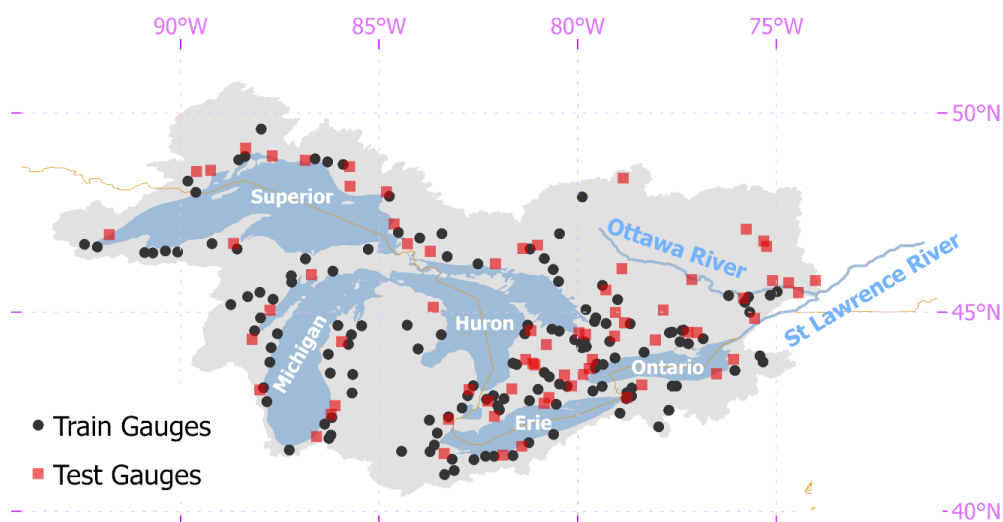
200

201    **2. Data**

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

203    in the St. Lawrence River basin (Figure 1). We note that this region is of similar spatial scale to other

204  benchmarking datasets for DL rainfall-runoff models (e.g., CAMELS-GB; Coxon et al., 2020). For direct

205  comparability to previous results from the GRIP-GL, all data for these watersheds are taken directly from

206  the work in Mai et al. (2022) and include daily streamflow time series, meteorological forcings, geophysical

207  attributes for each watershed, and auxiliary hydrologic fluxes. Daily streamflow were gathered from the

208  U.S. Geological Survey (USGS) and Water Survey Canada (WSC) between January 2001 and December

209  2017. All streamflow gauging stations have a drainage area greater than or equal to 200 km$^2$ and less than

210  5% missing data in the study period. The watersheds are evenly distributed across the five lake basins and

211  the Ottawa River basin, and they represent a range of land use/land cover types and degrees of hydrologic

212  alteration from human activity. In the experiments described further below, 141 of the watersheds are

213  designated as training sites, and the remaining 71 watersheds are used for testing (see Figure 1). In addition,

214  the period between January 2001 to December 2010 is reserved for model training (termed the training

215  period), and the period between January 2011 – December 2017 is used for model testing (termed the testing

216  period).

217



219  **Figure 1.** Great Lakes domain, with training and testing streamflow gauges used throughout this study.

220

221    Meteorological forcings are taken from the Regional Deterministic Reanalysis System v2 (RDRS-v2),

222    which is an hourly, 10 km dataset available across North America (Gasset et al., 2021). Hourly precipitation,

223    net incoming shortwave radiation ($R_s$), specific humidity (SH), surface pressure (SP), wind speed, and

224    temperature are aggregated into a basin-wide daily precipitation average, daily $R_s$ average, daily SH average,

225    daily SP average, daily wind speed average, and daily minimum and maximum temperature. We note that

226    the precipitation data from RDRS-v2 is produced from the Canadian Precipitation Analysis (CaPA), which

227    combines available surface observations of precipitation with a short-term reforecast provided by the 10

228    km Regional Deterministic Reforecast System. That is, the precipitation data is not model based, but rather

229    is based on gauged data and spatially interpolated using information from modeled output.

230

231    Geophysical attributes for each watershed were collected from a variety of sources. Basin-average statistics

232    of elevation and slope were derived from the HydroSHEDS dataset (Lehner et al., 2008), which provides a

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

234    gathered from the Global Soil Dataset for Earth System Models (GSDE; Shangguan et al., 2014), which is

235    available at a 30 arcsec resolution. Land cover data at a 30 m resolution and based on Landsat imagery from

236    2010-2011 were derived from the North American Land Change Monitoring System (NALCMS, 2017).

237    These geophysical datasets were used to derive basin-averaged attributes for each watershed, listed in Table

238    1.

239

240    **Table 1**. Watershed attributes used in the deep learning models developed in this work (adapted from Mai
241    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-broadleaf-forest** | Fraction of land covered by "Temperate-or-sub-polar-broadleaf-forest" |
| **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) |

242

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

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

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

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

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

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

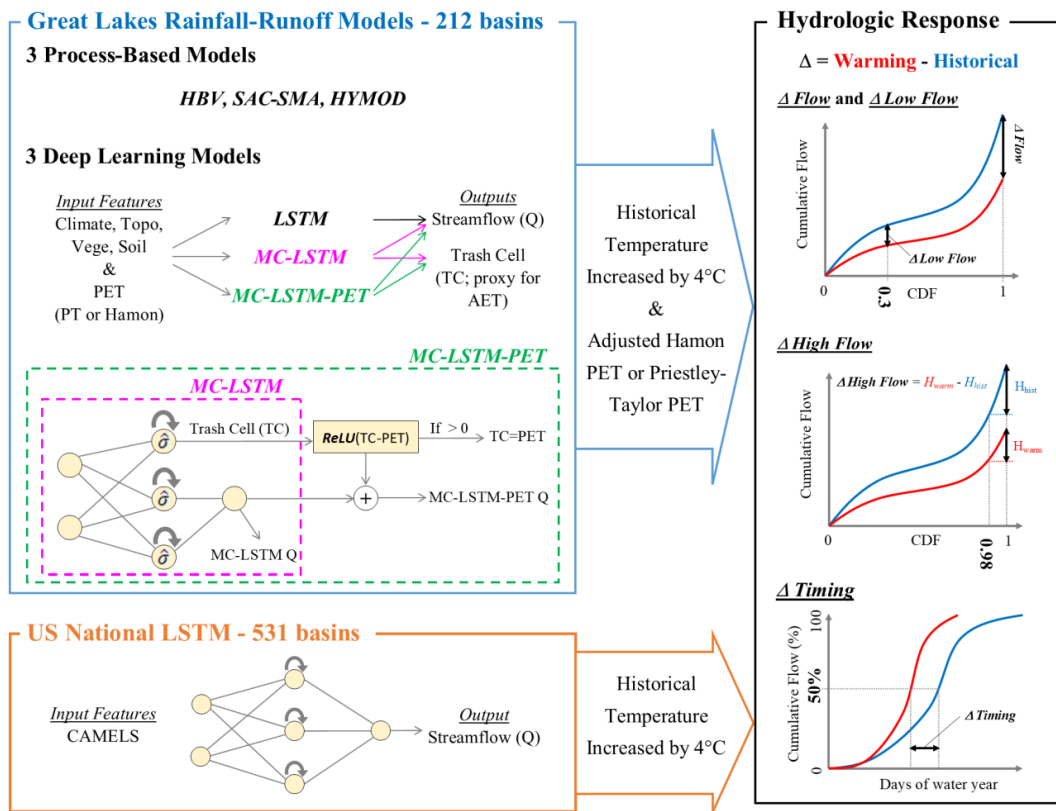249    2003-2017 and were aggregated to basin-wide totals for each watershed.

250

251    **3. Methods**

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

253    will overestimate hydrologic losses under warming because of an overreliance on historical correlations

254    between temperature and PET, while this effect will be lower in PIML-based rainfall-runoff models

255    designed to better account for water loss in the system. To conduct this experiment, we develop three

256    different DL rainfall-runoff models to predict daily streamflow across the Great Lakes region, as well as

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

258    based or temperature-based estimate of PET. The DL models include a regional LSTM very similar to the

259    model in Mai et al., (2022), an MC-LSTM that conserves mass, and a new variant of the MC-LSTM that

260    also respects the relationship between PET and water loss (termed MC-LSTM-PET). After comparing

261    historical model performance, we force all models with climate change scenarios composed of historical

262    precipitation and historical but warmed temperatures, as well as PET based on those warmed temperatures.

263    This is a similar approach to that taken in SW22, but in contrast to that study this work 1) focuses on the

264    magnitude of streamflow response to warming under two different PET formulations; 2) considers a

265    different set of physics-informed DL models in which the architecture (rather than the inputs or targets) of

266    the model are changed to better preserve physical plausibility under unprecedented climate change; and 3)

267    evaluates an expanded set of hydrologic metrics to better understand both the plausibility and the variability

268    of climate change responses across the different models. Finally, in a subset of the analysis, we also utilize

269    a fourth DL model, the LSTM used in SW22 that was previously fit to 531 basins across the contiguous

270    United States (Kratzert et al. 2021), which uses daily precipitation, maximum and minimum temperature,

271    radiation, and vapor pressure as input but not PET. This model is used to evaluate whether a DL model fit

272    to many more watersheds that span a more diverse gradient of climate conditions behaves differently under

273   warming than an LSTM fit only to locations in the Great Lakes basin. Figure 2 presents an overview of our

274   experimental design.

275



276

**Figure 2.** Overview of experiment design.

278

279   **3.1. Models**

280   **3.1.1. Benchmark Conceptual Models**

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

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

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

284   (Anderson, 1976). These models are developed as lumped, conceptual models for each watershed. We

285 calibrate the models with the genetic algorithm from Wang et al. (1991) to maximize the Kling-Gupta

286 Efficiency (KGE; Gupta et al. 2009), using a population size equal to 100 times the number of parameters,

287 evolved over 100 generations, and with a spin-up period of 1 year. Each benchmark model is calibrated

288 separately to each of the 141 training sites using the temporal train/test split described in Section 2.

289 Benchmark models are developed for the 71 testing sites in two ways: 1) separate models are trained for

290 the testing sites during the training period; and 2) each testing site is assigned a donor from among the 141

291 training sites, and the calibrated parameters from that donor site are transferred to the testing site. The first

292 of these approaches enables a comparison between DL models fit only to the training sites to benchmark

293 models developed for the testing sites, i.e., a spatial out-of-sample versus in-sample comparison. The

294 second of these approaches enables a more direct spatial out-of-sample comparison between DL and

295 benchmark models. We note that donor sites were used to assign model parameters to testing sites in the

296 benchmarking study of Mai et al. (2022), and to retain direct comparability to the results of that work we

297 use the same donor sites for each testing site. Donor sites were selected based on spatial proximity, while

298 also prioritizing donor sites that were nested within the watershed of the testing site.

299

300 **3.1.2. LSTM**

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

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

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

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

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

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

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

308 structure of the LSTM is given as follows:

309

14

310    $i[t] = \sigma(W_i x[t] + U_i h[t-1] + b_i)$                    (Eq. 1.1)

311    $f[t] = \sigma(W_f x[t] + U_f h[t-1] + b_f)$                    (Eq. 1.2)

312    $g[t] = tanh(W_g x[t] + U_g h[t-1] + b_g)$                    (Eq. 1.3)

313    $o[t] = \sigma(W_o x[t] + U_o h[t-1] + b_o)$                    (Eq. 1.4)

314    $c[t] = f[t] \odot c[t-1] + i[t] \odot g[t]$                    (Eq. 1.5)

315    $h[t] = o[t] \odot tanh(c[t])$                    (Eq. 1.6)

316    $y[T] = ReLU(W_y h[T] + b_y)$                    (Eq. 1.7)

317

318    Here, the input gate ($i[t]$) controls how candidate information ($g[t]$) from inputs and previous hidden states

319    flows to the current cell state ($c[t]$); the forget gate ($f[t]$) enables removal of information within the cell

320    state over time; and the output gate ($o[t]$) controls information flow from the current cell state to the hidden

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

322    streamflow predictions, $h[T]$ at the last time step in the sequence is passed through a fully connected layer

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

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

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

326    precipitation to be conserved within this architecture.

327

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

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

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

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

332    steps in the input sequences $x$. All input features are standardized before training (by subtracting the mean

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

334    do not standardize the observed streamflow, besides dividing my drainage area to represent streamflow in

335    units of millimeters.

336

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

338    during the training period:

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

340    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

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

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

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

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

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

346    the training sites. Network weights are tuned using the ADAM optimizer (Kingma & Ba, 2015). The model

347    is trained 10 separate times with different random initializations to account for uncertainty in the training

348    process.

349

350    For the evaluation of streamflow projections under climate change, we also use an LSTM taken from

351    Kratzert et al. (2021) and employed in SW22, which was fit to 531 basins across the contiguous United

352    States (hereafter called the National LSTM). This model was trained using a different set of data compared

353    to our Great Lakes LSTM but also used a mix of dynamic and static features, all of which were drawn from

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

355    2015). This model uses daily precipitation, maximum and minimum temperature, shortwave downward

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

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

358

### 3.1.3. MC-LSTM

Following Hoedt et al. (2021) and Frame et al. (2021b), we adapt the architecture of the LSTM into a mass conserving MC-LSTM that preserves the water balance within the model, i.e., the total quantity of precipitation entering the model is tracked and redistributed to streamflow and losses from the watershed. Using similar notation as for the LSTM above, the model structure is given as follows:

$$\boldsymbol{i}[t] = \hat{\sigma}(\boldsymbol{W_i}\boldsymbol{x}[t] + \boldsymbol{U_i}\boldsymbol{c}[t-1] + \boldsymbol{V_i}\boldsymbol{a}[t] + \boldsymbol{b_i}) \tag{Eq. 3.1}$$

$$\boldsymbol{o}[t] = \sigma(\boldsymbol{W_o}\boldsymbol{x}[t] + \boldsymbol{U_o}\boldsymbol{c}[t-1] + \boldsymbol{V_o}\boldsymbol{a}[t] + \boldsymbol{b_o}) \tag{Eq. 3.2}$$

$$\boldsymbol{R}[t] = \hat{\sigma}(\boldsymbol{W_R}\boldsymbol{x}[t] + \boldsymbol{U_R}\boldsymbol{c}[t-1] + \boldsymbol{V_R}\boldsymbol{a}[t] + \boldsymbol{b_R}) \tag{Eq. 3.3}$$

$$\boldsymbol{m}[t] = \boldsymbol{R}[t]\boldsymbol{c}[t-1] + \boldsymbol{i}[t]\boldsymbol{x}[t] \tag{Eq. 3.4}$$

$$\boldsymbol{c}[t] = (1 - \boldsymbol{o}[t]) \odot \boldsymbol{m}[t] \tag{Eq. 3.5}$$

$$\boldsymbol{h}[t] = \boldsymbol{o}[t] \odot \boldsymbol{m}[t] \tag{Eq. 3.6}$$

Here, the inputs to the model are split between quantities **x**[t] to be conserved (i.e., precipitation), and non-conservative inputs **a**[t] (i.e., temperature, wind speeds, PET, catchment properties, etc.). Water in the system is stored in the *D*-dimensional vector **m**[t] and is updated at each time step based on water left over from the previous time step (**c**[t-1]) and water entering the system at the current time step (**x**[t]). The input gate **i**[t] and a redistribution matrix **R**[t] are designed to ensure water is conserved from $\boldsymbol{c}[t-1]$ and $\boldsymbol{x}[t]$ to **m**[t], by basing these quantities on a normalized sigmoid activation function that ensures a column-normalized **R**[t] and **i**[t] summing to unity.

The mass in $\boldsymbol{m}[t]$, which is stored across *D* elements in the vector, is then distributed to the output of the hidden layer, $\boldsymbol{h}[t]$, or the next cell state, $\boldsymbol{c}[t]$. To account for water losses from evapotranspiration or other sinks, one element of the *D*-dimensional vector $\boldsymbol{h}[t]$ is considered a 'trash cell', and the output of this cell

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

384    water mass:

385

386    $y[T] = \sum_{d=1}^{D-1} h_d[T]$           (Eq. 4)

387

388    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 $t=1,..,T$

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

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

391    with random initialization).

392

393    **3.1.4. MC-LSTM-PET**

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

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

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

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

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

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

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

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

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

403

404    $y[t] = \sum_{d=1}^{D-1} h_d[t] + ReLU(h_D[t] - PET[t])$           (Eq. 5)

405

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

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

18

408 LSTM (Eq. 4) if water in the trash cell is less than PET. This approach assumes that the maximum allowable

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

410 deep groundwater percolation that remains disconnected from the stream; lateral groundwater flows out of

411 the watershed; human diversions). However, given that evapotranspiration accounts for the vast majority

412 of water lost in most hydrologic systems, this assumption is likely reasonable in most cases. The MC-

413 LSTM-PET was trained in the same way as the LSTM (i.e., same inputs, loss function, training and test

414 sets, hyperparameter selection process, number of ensemble members with random initialization).

415

416 **3.2. Model Performance Evaluation**

417 As noted previously, 141 of the watersheds are designated as training sites, and the remaining 71 watersheds

418 are used for testing. In addition, the training and testing periods were restricted to January 2001 -December

419 2010 and January 2011 – December 2017, respectively. This provides three separate ways to evaluate model

420 performance:

421 • Temporal validation - Performance across models is evaluated at training sites during the testing

422 period.

423 • Spatial validation - Performance across models is evaluated at testing sites during the training

424 period.

425 • Spatiotemporal validation - Performance across models is evaluated at testing sites during the

426 testing period.

427

428 All three evaluation strategies are utilized. For benchmark process-based models that are calibrated locally

429 on a site-by-site basis, we consider model versions that are transferred to testing sites from training sites,

430 as well as models that are trained to the testing sites directly (see Section 3.1.1). The former can be used

431 for all three evaluation strategies above, while the latter can only be used for temporal validation at the

432 testing sites.

433

434 Several metrics are considered for model evaluation, including percent bias (PBIAS), the Nash-Sutcliffe

435 Efficiency (NSE; Nash and Sutcliffe, 1970), Kling-Gupta Efficiency (KGE; Gupta et al. 2009), top 2%

436 peak flow bias (FHV; Yilmaz et al. 2008), and bottom 30% low flow bias (FLV; Yilmaz et al. 2008). Each

437 metric is calculated separately for training and testing periods for each site. For the DL models, all results

438 are estimated from the ensemble mean from 10 separate training trials.

439

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

441 observations of AET from the GLEAM database. We note that AET data were not used to train any of the

442 models. For the process models, AET is a direct output of the model and so can immediately be extracted

443 for comparison, but AET is not directly simulated by the MC-LSTM or MC-LSTM-PET. Instead, we

444 assume water delegated to the trash cell permanently leaves the system because of evapotranspiration.

445 Several metrics are used to compare model-based AET to GLEAM AET, including KGE, correlation, and

446 PBIAS, and the comparison is conducted for training sites during the training period and under temporal,

447 spatial, and spatiotemporal validation (as described above). Similar to streamflow, all AET results for the

448 MC-LSTM and MC-LSTM-PET are based on the ensemble mean of water delegated to the trash cell from

449 the 10 separate training trials.

450

451 **3.3. Evaluating Hydrologic Response under Warming**

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

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

454 budget-based approach; Priestley and Taylor, 1972). PET (in mm/day) under the Hamon method is

455 calculated as follows (Shaw and Riha, 2011):

456

457 $PET_H = \alpha_H \times 29.8 \times Hr_{day} \frac{e_{sat}}{T_a + 273.2}$ (Eq. 6)

458     $e_{sat} = 0.611 \times exp\left(\frac{17.27 \times T_a}{237.3 + T_a}\right)$ (Eq. 7)

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

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

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

462

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

464

465     $PET_{PT} = \alpha_{PT}\left(\frac{\Delta(T_a) \times (R_n - G)}{\lambda(\Delta(T_a) + \gamma)}\right) \times 1000$ (Eq. 8)

466

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

468     $T_a$, $\gamma$ is the psychrometric constant (kPa/°C), $\lambda$ is the volumetric latent heat of vaporization (MJ/m$^3$), $R_n$ is

469     the net radiation (MJ/m$^2$-day) equal to the difference between net incoming shortwave ($R_{ns}$) and net

470     outgoing longwave ($R_{nl}$) radiation, $G$ is the heat flux to the ground (MJ/m$^2$-day), and $\alpha_{PT}$ is a dimensionless

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

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

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

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

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

476     are derived from the RDRS-v2 (see Section 2). We note that using $\alpha_H = 1.2$ and $\alpha_{PT} = 1.1$ leads to very

477     similar PET estimates between the Hamon and Priestley-Taylor methods under baseline climate conditions,

478     helping to ensure their comparability.

479

480     We then develop a simple climate change scenario in which the historical minimum and maximum

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

21

482    warmed temperatures. We focus the climate change assessment on training period data at the training sites,

483    so that any differences in climate change projections that emerge between the DL and process models are

484    due to model structural differences and not the effects of spatiotemporal regionalization. In the Priestly-

485    Taylor method, we maintain historical values for $R_s$ to isolate how changes in temperature and its effect on

486    $\Delta(T_a)$ and $R_{nl}$ influence changes in PET. The use of historical $R_s$ is supported by the results from CMIP5

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

488

489    We also develop a similar climate change scenario for the National LSTM, which uses five dynamic input

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

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

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

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

494    2015). Thus, to run the National LSTM under warming, we also adjust the vapor pressure input based on

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

496

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

498    warming scenarios above. We also consider changes to some of the static input features that depend on

499    temperature and PET (e.g., pet_mean, aridity, t_mean, frac_snow; see Table 1) and run all models using

500    two settings: 1) with climate changes only to the dynamic features, and 2) with climate changes to both

501    dynamic and static features.

502

503    Ultimately, for each model we compare hydrologic projections under the warmed scenario to their values

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

505    CAMELS dataset within the Great Lakes Basin. We examine four different metrics for this comparison,

506    including:

22

507    • AVG.Q: the average runoff across the entire series.

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

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

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

511    the year by which half of the total annual flow has passed.

512

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

514    projections but similar warming while process-based and PIML models can, we would expect that under

515    the LSTM using both PET projections, average flow will decline significantly and with similar magnitude

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

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

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

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

520    process model projections more closely across the two different PET projections, at least in terms of the

521    difference in magnitude of average streamflow declines. For comparison, we also explore the differences

522    in low flow (FLV), high flow (FHV), and timing (COM) metrics across all model versions, where we have

523    less reason to anticipate how DL and process models will differ in their projections and across PET

524    formulations.

525

526    **4. Results**

527    **4.1. Model Performance Evaluation**

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

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

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

531    Priestley-Taylor PET, but there is little difference in performance for the models fit with Hamon PET (see
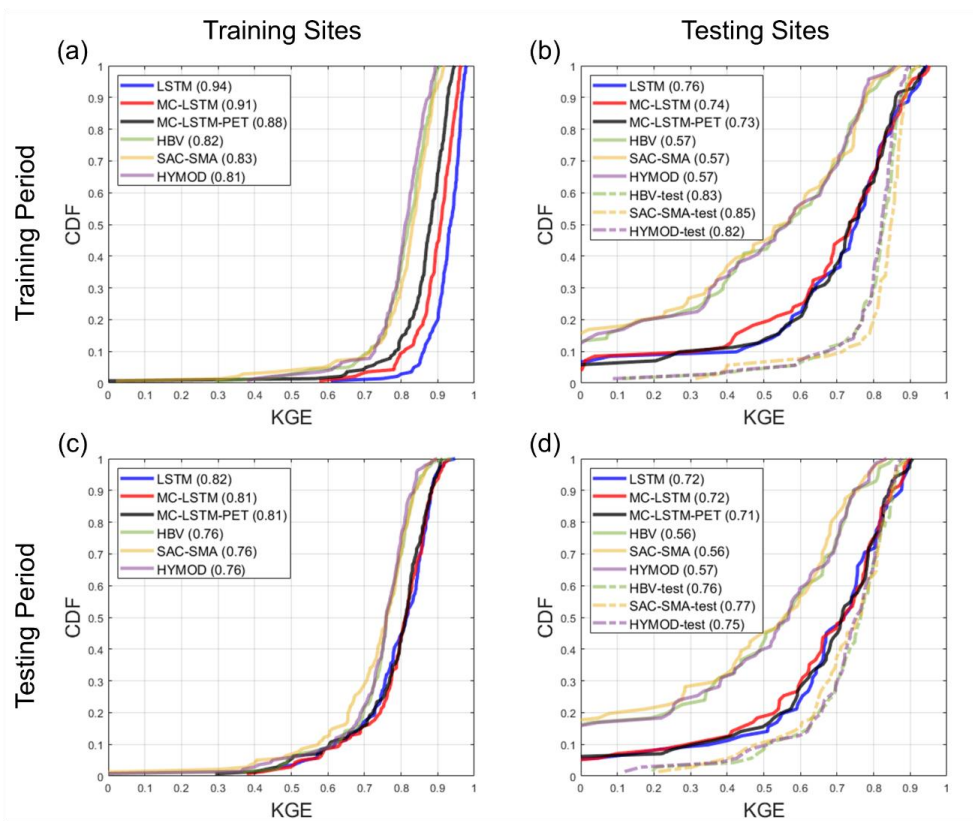
532    Figure S1). For the process-based models, we show results for models fit to the training sites and then used

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

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

535    to the testing sites.

536

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

538    perform very well (Figure 3a). Across the three process models, the median KGE is 0.82, 0.83, and 0.81

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

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

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

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

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

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

545    between 0.56-0.57 when streamflow at the testing sites is estimated with donor models from nearby gauged

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

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

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

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

550    performance of the DL models. We even see that for approximately 10% of testing sites during the training

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

552

**Figure 3.** The distribution of Kling-Gupta efficiency (KGE) for streamflow estimates across sites from each model at the (a) 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 2 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, with the exception of PBIAS during the training period. 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 KGE, NSE, and FLV, 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. We note that percent biases for FLV are high because the absolute magnitude of low flows is small, so small absolute biases still lead to large percent biases.

25

568

569 **Table 2.** The median KGE, NSE, PBIAS, FHV, and FLV for streamflow across testing sites for the training
570 and testing periods for all models (excluding the process models fit to the testing sites). The metric from
571 the best performing model in each period is bolded. All models are trained using Priestley-Taylor PET.
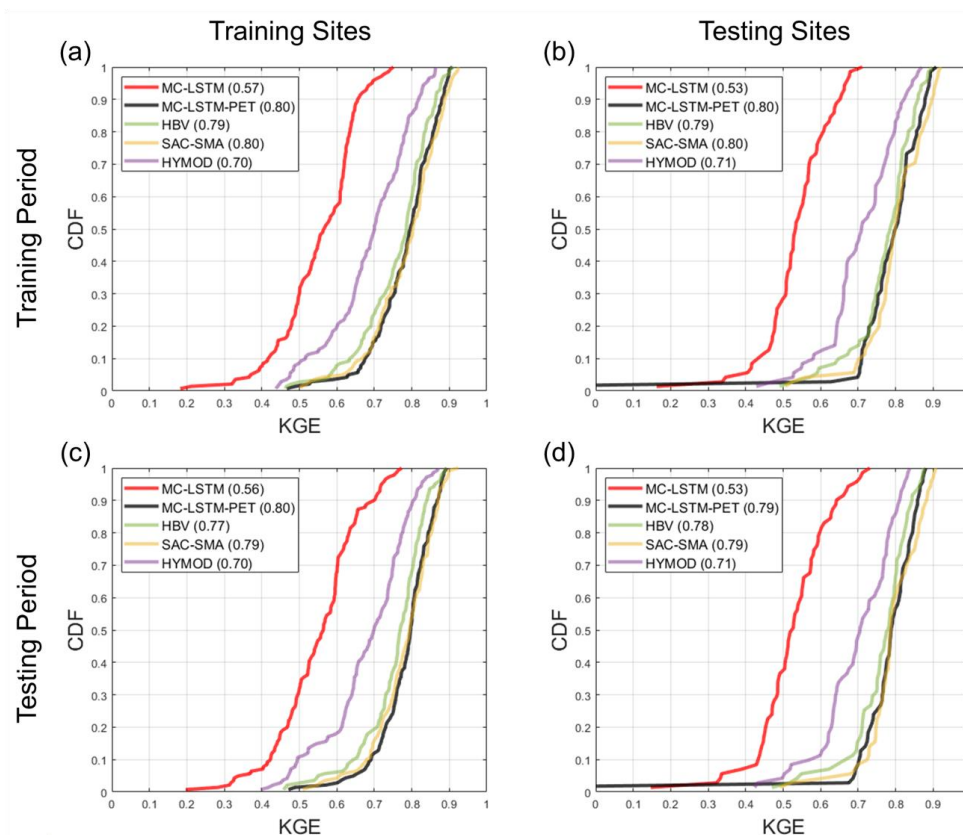
| Model | Testing Sites: Training Period | | | | | Testing Sites: Testing Period | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | KGE | NSE | PBIAS | FHV | FLV | KGE | NSE | PBIAS | FHV | FLV |
| LSTM | **0.76** | **0.77** | 9.66 | 17.58 | **30.98** | **0.72** | **0.68** | 12.15 | 26.01 | **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 | 0.57 | 0.42 | **8.41** | 32.61 | 50.41 | 0.56 | 0.45 | 11.24 | 36.29 | 46.67 |
| SAC-SMA | 0.57 | 0.43 | 11.03 | 34.54 | 42.08 | 0.56 | 0.41 | 12.13 | 36.74 | 49.29 |
| HYMOD | 0.57 | 0.41 | 9.58 | 32.70 | 52.24 | 0.57 | 0.45 | 11.16 | 36.34 | 53.62 |

572

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

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

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

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

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

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

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

580    0.77-0.80). Performance is slightly worse for HYMOD. Notably, the MC-LSTM-PET exhibits very similar,

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

582    site. In contrast, the MC-LSTM performs the worst of all models, with median KGE values ranging between
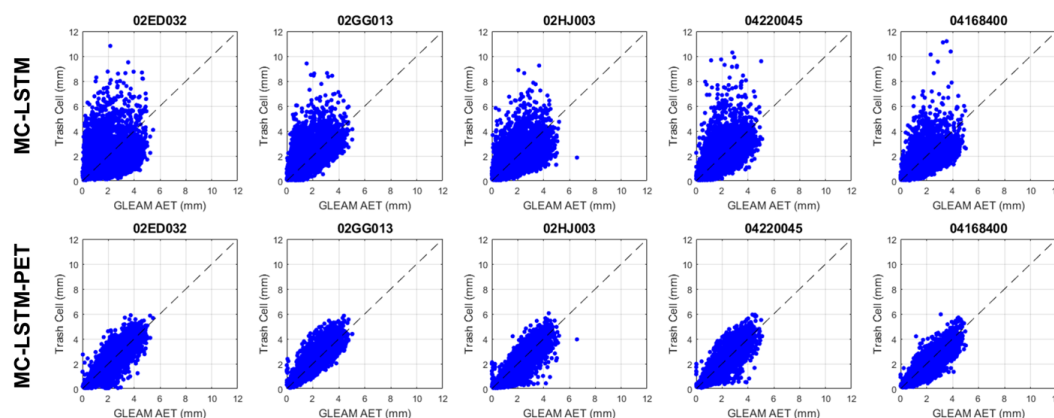
583    0.53-0.57.

584

585

**Figure 4.** The Kling-Gupta efficiency (KGE) for AET estimated from each model at the (a) 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 observed AET versus water allocations to the trash cell for the two models from five randomly sampled testing sites across both training and testing periods. Trash cell water from the MC-LSTM is not only more scattered around observed AET compared to the MC-LSTM-PET, but it also exhibits many outlier values that are two to five times larger than observed AET. The MC-LSTM-PET follows the variability of AET much more closely, with virtually no outliers

599    that exceed AET by large margins. This suggests that the PET constraint on the trash cell in the MC-LSTM-

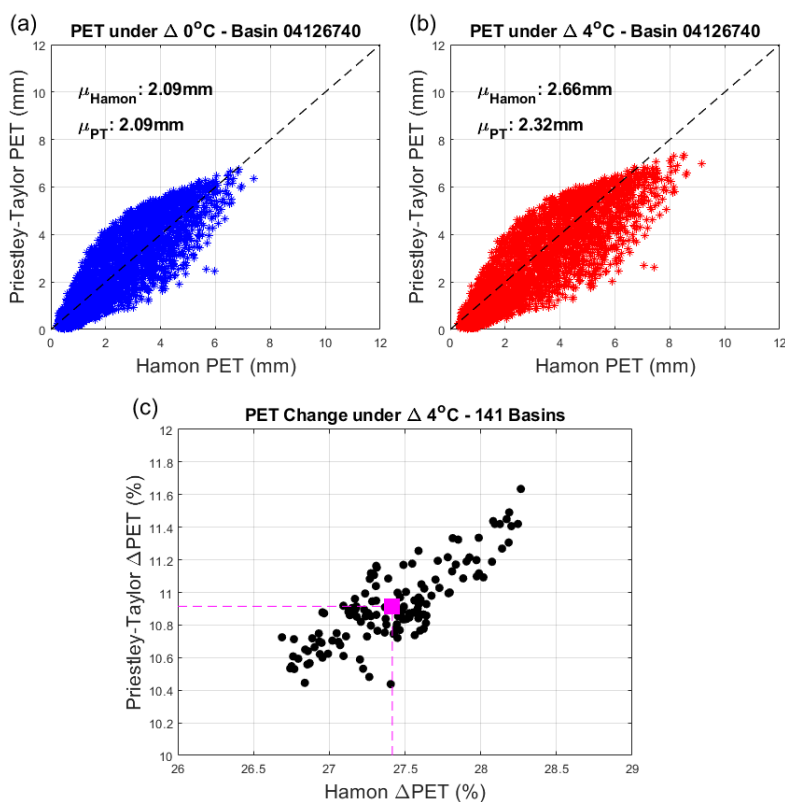600    PET helps water allocated to that cell more faithfully represent an ET sink in the DL model.



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

## 4.2. Evaluating Hydrologic Response under Warming

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

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

609    model structure and not spatiotemporal regionalization. First, we show the differences in historic and

610    projected PET when using the Hamon and Priestley-Taylor methods (Figure 6). For the training period

611    without any temperature change, PET estimated from the two methods is very similar (shown at one sample

612    location for demonstration; Figure 6a). However, under the scenario with 4 °C of warming, Hamon-based

613    PET is significantly larger than Priestley-Taylor based PET (Figure 6b). On average, this difference reaches

614    ~16% across all training sites and exhibits very little variability across locations (Figure 6c). The primary

615    reason for the difference in projected change in PET is that the Hamon method attributes PET entirely to

616    temperature, while only a portion of PET is based on temperature in the Priestley-Taylor method, with the

617    rest based on $R_n$. It is worthwhile to note that $R_n$ does change with temperature through its effects on net

618    outgoing longwave radiation, but these changes are small.

619



620

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

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

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

628    projections using Hamon and Priestley-Taylor PET, respectively, while the rows of Figure 7 show the

629    distribution of changes (as a percentage) in different streamflow attributes (AVG.Q, FLV, FHV, COM)

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

631    warming, while Figure S4 show the same results when both the dynamic and the static climate properties

632    are updated with warming.

633

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

635     exhibit similar changes in average streamflow to the process-based models, with the median ΔAVG.Q

636     across sites ranging between -23% and -17% across all models. However, when using Priestley-Taylor PET,

637     larger differences in the distribution of ΔAVG.Q emerge. Across all three process models, the median

638     ΔAVG.Q is between -9% to -5%, and very few locations exhibit ΔAVG.Q less than -20%. Conversely, the

639     LSTM shows a median water loss of -20% under Priestley-Taylor PET and a very similar distribution of

640     water losses regardless of whether Hamon or Priestley-Taylor PET was used. The MC-LSTM is also

641     relatively insensitive to PET, and as compared to the process models, the MC-LSTM tends to predict

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

643     the MC-LSTM-PET model achieves water loss that is significantly smaller under Priestley-Taylor PET

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

645

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

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

648     significantly on changes to FLV and bound the changes predicted by the DL models. HBV and HYMOD

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

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

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

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

653     and so small changes in millimeters of flow lead to large percent changes.

654

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

656     streamflow timing (COM; Figure 7g,h) are relatively consist, with the process models exhibiting larger

657     declines in high flows and earlier shifts in streamflow timing compared to the DL models. The choice of
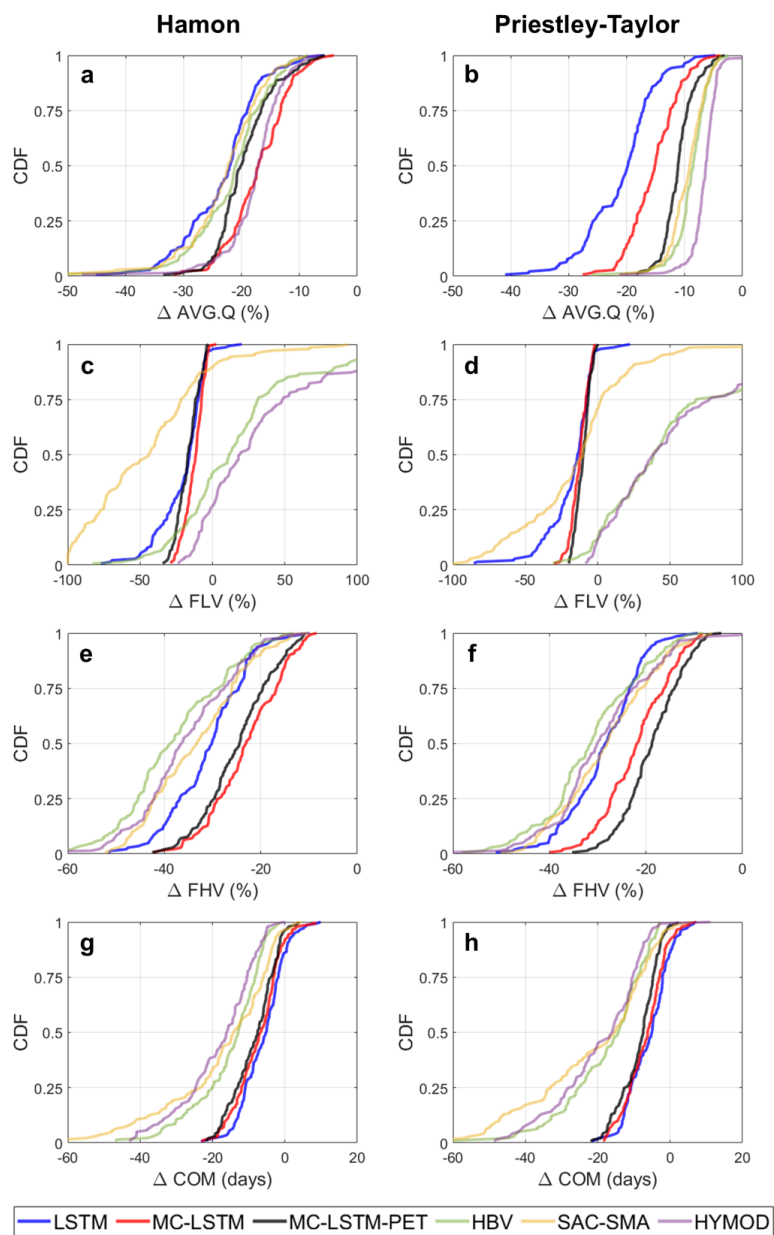
658    PET method has a moderate impact on process-model based changes in FHV, with larger declines under

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

660    although the LSTM predicts changes in FHV closest to the process models. For COM, the process models

661    show a wide range of variability in projected change across sites, from no change to 60 days earlier. For

662    the DL models the range of change is much narrower, and the median change in COM is almost a week

663    less than the median change across the process models. The method of PET estimation has relatively little

664    impact on both process model and DL based estimates of change in COM.

665

666    We note that if the static watershed properties (pet_mean, aridity, t_mean, frac_snow; see Table 1) are also

667    changed to reflect warmer temperatures and higher PET, all three DL models exhibit unrealistic water gains

668    for between 15%-40% of locations depending on the model and PET method, with the most water gains

669    occurring under the LSTM (Figure S4). These results suggest that changing the static watershed properties

670    associated with long-term climate characteristics can degrade the quality of the projections, at least when

671    the climate changes are large and the range of average temperature and PET in the training set is limited.

672    We also note that the results in Figure 7 are largely unchanged if based on projections for testing sites in

673    the testing period (Figure S5).

674

**Figure 7.** The distribution of change in (a,b) AVG.Q, (c,d) FLV, (e,f) FHV, and (g,h) 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 DL models, changes were only made to the dynamic inputs (i.e., no changes to static inputs).

681    One reason why the Great Lakes LSTM exhibits excessive hydrologic losses under warming could be that

682    the model was trained using sites that are confined to a limited range of temperature and PET values found

683    in the Great Lakes basin (spanning approximately 40.5°-50°N), and so is ill-suited to extrapolate hydrologic

684    response under warming conditions that extend beyond this range. To evaluate this hypothesis, we examine

685    changes to AVG.Q, FLV, FHV, and COM under 4°C warming at the 29 CAMELS watersheds within the

686    Great Lakes basin using the National LSTM (Figure 8). For comparison, we also examine similar changes

687    under all six Great Lakes DL and process models at 17 of those 29 CAMELS basins that were used in the

688    training and testing sets for the Great Lakes models, and also separate out the National LSTM projections

689    for those 17 sites. Note that in Figure 8, the National LSTM projections do not differ between Hamon and

690    Priestley Taylor PET, because PET is not an input to that model.
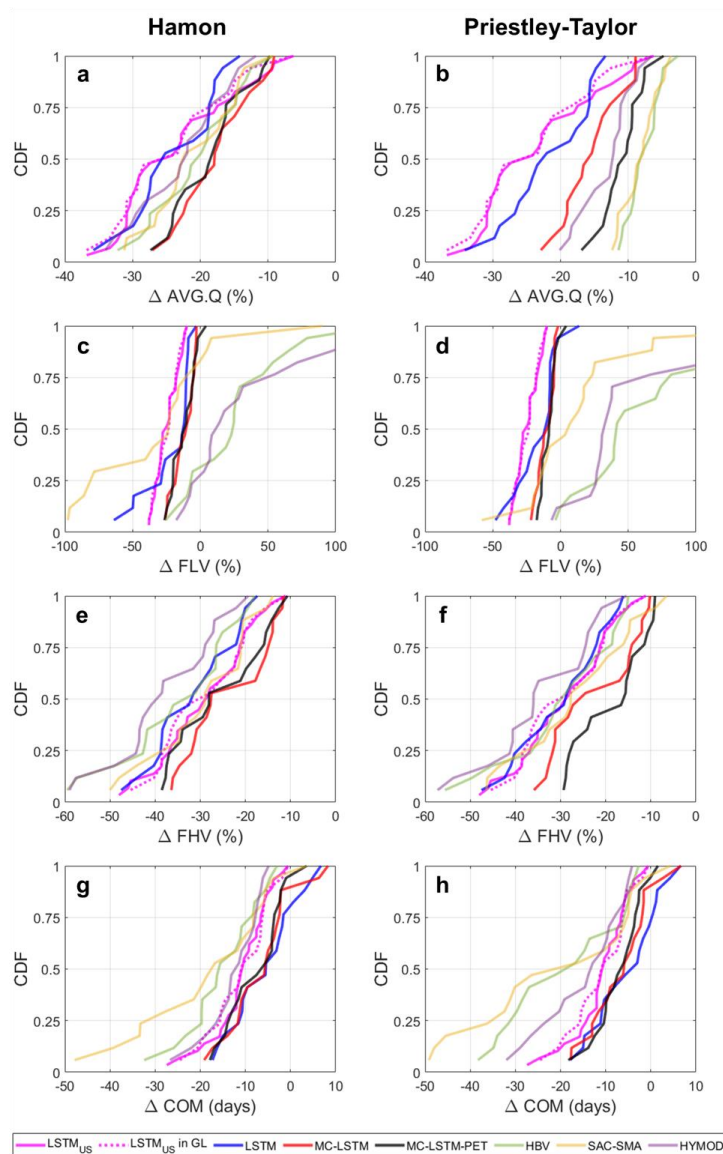
691

692    The National LSTM was trained to watersheds across the contiguous United States (spanning

693    approximately 26°-49°N), and so was exposed to watersheds with much warmer conditions and higher PET

694    during training. However, we find that the National LSTM still projects very large declines in AVG.Q. For

695    the 29 CAMELS watersheds in the Great Lakes basin, the median decline in AVG.Q under the National

696    LSTM is approximately 25%, which is moderately larger than the median projections of loss under the

697    process models using Hamon PET and much larger than the process model losses under Priestley-Taylor

698    PET (Figure 8a,b). We also see larger declines in FLV under the National LSTM as compared to the other

699    Great Lakes DL models (Figure 8c,d). The National LSTM projects changes in FHV (Figure 8e,f) and COM

700    (Figure 8g,h) that are similar to the process models, and for COM, the projections are closer to the process

701    models than for any Great Lakes DL model. In addition, the hydrologic projections are stable under the

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

703    under warming (see Figure S6), in contrast to the Great Lakes DL models. Therefore, the use of more

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

705    inputs are shifted significantly to reflect new climate conditions. However, as shown in Figure 8a,b, this

706    benefit does not mitigate the tendency for the National LSTM to overestimate water loss under warming.
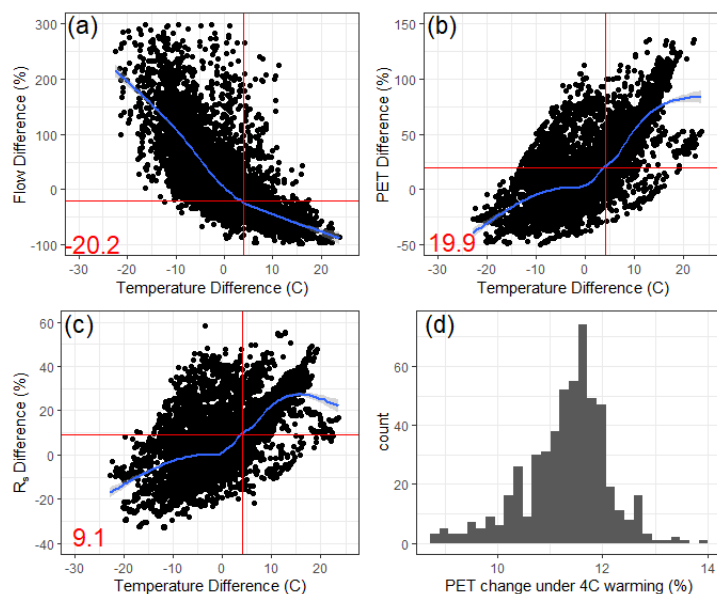
707



708

**Figure 8.** The distribution of change in (a,b) AVG.Q, (c,d) FLV, (e,f) FHV, and (g,h) 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 DL 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

713  models only, results differ when using (a,c,e,f) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the
714  National LSTM, changes were made only to the dynamic inputs.
715

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

717  to examine how average streamflow, (Priestly-Taylor estimated) PET, and $R_s$ vary across all 531 CAMELS

718  watersheds of different average temperatures, and compare this variability to projected changes in PET at

719  each site under warming. Specifically, we compare the difference in long-term (1980-2014) average

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

721  dataset with average long-term precipitation within 1% of each other, and plot these differences against the

722  differences in average temperature across each pair. The results show that the difference in average

723  streamflow across watersheds with similar precipitation becomes negative when the difference in

724  temperature is positive (i.e., warmer watersheds have less flow on average), and that when the difference

725  in average temperature reaches 4°C, flows differ by about 20% on average (Figure 9a). This is very similar

726  to the projected median decline in average streamflow seen for the National LSTM in Figure 8. We also

727  note that average PET increases by approximately 20% between watersheds that differ in average

728  temperature by 4°C (Figure 9b). However, higher PET in warmer watersheds is related both to the direct

729  effect of temperature on vapor pressure deficit, as well as to the fact that higher incoming solar radiation

730  co-occurs in warmer watersheds ($R_s$ is approximately 9% higher across watershed pairs that differ by 4°C;

731  Figure 9c). Using the Priestley-Taylor method, we estimate that average PET would only increase by

732  between 9-14% (median of 11.5%) if temperatures warm by 4°C and $R_s$ is held at historic values, while $R_n$

733  is increased slightly due to declines in net outgoing longwave radiation with warming (Figure 9d). However,

734  the National LSTM appears to convolute the effects of temperature and $R_s$ and cannot separate out their

735  effects on ET-based water loss, leading to larger projected streamflow losses under 4°C warming than

736  changes in PET would warrant. This is possibly because of the very strong correlation between at-site daily

737  temperature and $R_s$ historically (median correlation of 0.85 across all CAMELS watersheds).

738

**Figure 9.** The percent difference in long-term (1980-2014) average (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 an analysis that evaluates the physical plausibility of future streamflow projections under climate change using DL rainfall-runoff models. The basis for this evaluation is anchored to the assumption that differences in streamflow projections should emerge under very different projections of future PET, and that realistic projections of future 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 future streamflow projections should be able to respond to lower energy-budget based PET projections under warming and, all else equal, project 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 projections of future PET under warming. This

36

758    discrepancy in future projections emerged despite the fact that the standard LSTM was a far better model

759    for streamflow estimation in ungauged basins compared to three process-based models under historic

760    climate conditions. In addition, the National LSTM trained to a much larger set of watersheds (531 basins

761    across 23° of latitude) using temperature, vapor pressure, and $R_s$ directly (rather than PET) also estimated

762    water loss under warming that far exceeded the losses estimated with process models forced with energy

763    budget-based PET. Since water losses estimated using energy budget-based PET are generally considered

764    more realistic (Lofgren et al., 2011; Shaw and Riha, 2011; Lofgren and Rouhana, 2016; Milly and Dunne,

765    2017; Lemaitre-Basset et al. 2022), this result casts doubt over the physical plausibility of the LSTM

766    projection.

767

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

769    streamflow responses under climate change while still maintaining superior prediction skill compared to

770    process models, at least in some cases. In particular, a mass conserving LSTM that also respected the limits

771    of water loss due to ET (the MC-LSTM-PET) was able to project changes in average streamflow that much

772    more closely aligned with process-model based estimates, while also providing competitive out-of-sample

773    performance across all models considered (including the other DL models). A more conventional MC-

774    LSTM that did not limit water losses by PET was less consistent with process-based estimates of change in

775    average streamflow. These results highlight the potential for PIML-based DL models to help achieve similar

776    performance improvements over process-based models as documented in recent work on DL rainfall-runoff

777    models (Kratzert et al., 2019a,b; Feng et al., 2020; Nearing et al., 2021) while also producing projections

778    under climate change that are more consistent with theory than non-PIML DL models.

779

780    An interesting result from this study was the disagreement in the change in high flows and streamflow

781    timing between all Great Lakes DL models and process models, the latter which estimated greater

782    reductions in high flows and larger shifts of water towards earlier in the year. Projections from the Great

783    Lakes DL models were also unstable if static climate properties of each watershed were changed under

37

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

785   changes to high flows and streamflow timing that were more like the process models than projections from

786   the Great Lakes DL models. While it is challenging to know which set of projections are correct for these

787   streamflow properties, these result overall favor projections from the National LSTM and highlight the

788   benefits of DL rainfall-runoff models trained to a larger set of diverse watersheds for climate change

789   analysis.

790

791   The MC-LSTM-PET model proposed in this work represents one (relatively simple) PIML-based

792   architectural change to an existing DL model in the hydrologic literature that can help better capture

793   physical constraints on water loss from hydrologic systems. However, other possibilities exist. For example,

794   the hard constraint in the MC-LSTM-PET could instead be imposed as a soft constraint through adjustments

795   to the loss function, where water losses in the trash cell that exceed PET are penalized. The MC-LSTM-

796   PET model could also be adjusted further to allow additional water losses in the trash cell related to human

797   water extractions from the watershed or other terminal sinks. A different approach would be to use learnable,

798   differentiable, process-based models with embedded neural networks (Jiang et al., 2020; Feng et al., 2022;

799   Feng et al., 2023), which can achieve similar performance to LSTMs but can also represent and output

800   different internal hydrologic fluxes. Further work is needed to evaluate the benefits and drawbacks of these

801   different PIML-based approaches, preferably on large benchmarking datasets such as CAMELS.

802

803   One important limitation of this study is how we constructed the climate change scenarios, with 4°C

804   warming but no change to net incoming shortwave radiation and slight decreases in net outgoing longwave

805   radiation with warming (i.e., slight increases in $R_n$). We did not consider any changes in net incoming

806   shortwave radiation because there is significant uncertainty in this term at local scales and its relationship

807   to local temperature change. Projections of net incoming shortwave radiation are highly variable across

808   space and can even differ in the direction of change, largely because of uncertainty in the representation of

809   clouds in climate models, future projections of aerosols, and the representation of cloud-aerosol interactions

810    (Chen, 2021; Coppola et al., 2021; Taranu et al., 2023). The relationship between local net radiation change

811    and local temperature change further depends on horizontal energy transport from other regions (Nordling

812    et al., 2021). In addition, the approximation we used for changes to net outgoing longwave radiation was

813    not designed to resolve all land-atmosphere energy balance feedbacks with changing atmospheric

814    composition under climate change. These uncertainties, along with uncertainties in energy-budget based

815    methods used to estimate PET (Greve et al. 2019; Liu et al., 2022), complicate future projections of

816    atmospheric drying power under warming. Regardless, the main finding of this work remains, namely that

817    DL models struggle to propagate different hypotheses of future PET scenarios into hydrologic projections

818    unless explicitly directed to do so.

819

820    Finally, we note that the results of this study do not entirely preclude the possibility that a standard LSTM,

821    fit to a sufficiently large set of diverse watersheds, could ultimately learn more physically realistic

822    projections under climate change. Our results with the National LSTM suggest that the signals between

823    temperature change and $R_s$ on water loss may be entangled, making it difficult for the model to estimate the

824    individual effects of changes to one of those terms (temperature) on water loss. However, it is possible that

825    the model would produce hydrologic projections that were more in line with theory if it was given 1) high

826    quality data on all terms related to water loss; and 2) future projections of these terms that were co-

827    developed in physically consistent ways (e.g., from physical climate models). The $R_s$ used in the National

828    LSTM was based on reanalysis and so may have had meaningful errors that drove the model to attribute

829    more water loss to warmer temperatures, and the scenario of warming given to the National LSTM (4°C

830    warming with no change in $R_s$) may violate the physical relationship between temperatures and $R_s$. While

831    outside the scope of the present study, we argue more work is needed to further explore the physical

832    plausibility of hydrologic projections with more standard LSTMs, with greater attention paid to the

833    meteorologic inputs used in the model under historical and future climate conditions.

834

835 **Acknowledgements**

837

838 **Competing Interests**

839 The authors declare no competing interests.

840

841 **Data and Code Availability Statement**

842 The code used for this project is available at https://doi.org/10.5281/zenodo.8190287. All data used to

843 train and evaluate the models are available at https://www.hydrohub.org/mips_introduction.html#grip-gl.

844

845 **References**

846 Allen, R.G., Pereira, L.S., Raes, D., et al. (1998) Crop Evapotranspiration-Guidelines for Computing
847 Crop Water Requirements-FAO Irrigation and Drainage Paper 56. FAO, Rome, 300(9): D05109.
848
849 Anderson, E. A. (1976). A point energy and mass balance model of a snow cover (NOAA Technical
850 Report NWS 19). Silver Spring, MD: National Oceanic and Atmosphere Administration.
851
852 Bergström, S. & Forsman, A. (1973) Development of a conceptual deterministic rainfall-runoff model.
853 Nordic Hydrol. 4, 147–170.
854
855 Beven, K. (2023). Benchmarking hydrological models for an uncertain future. Hydrological
856 Processes, 37( 5), e14882. https://doi.org/10.1002/hyp.14882
857
858 Boyle, D. P. (2001). Multicriteria calibration of hydrologic models, (Doctoral dissertation). Retrieved from
859 UA Campus Repository (http://hdl.handle.net/10150/290657), Tucson, AZ: The University of Arizona.
860
861 Burnash, R. J. (1995). The NWS river forecast system - catchment modeling. In Singh, V. (Ed.), Computer
862 Models of Watershed Hydrology (pp. 311-366). Littleton, CO: Water Resources Publication.
863
864 Campbell, M., Cooper, M. J. P., Friedman, K., & Anderson, W. P. (2015). The economy as a driver of
865 change in the Great Lakes - St. Lawrence basin. *Journal of Great Lakes Research*, *41*, 69–83.
866
867 Cayan, D. R., Kammerdiener, S. A., Dettinger, M. D., Caprio, J. M., & Peterson, D. H. (2001). Changes
868 in the Onset of Spring in the Western United States, *Bulletin of the American Meteorological*
869 *Society*, 82(3), 399-416. https://doi.org/10.1175/1520-0477(2001)082<0399:CITOOS>2.3.CO;2
870
871 Chen, L. Uncertainties in solar radiation assessment in the United States using climate models. Clim
872 Dyn 56, 665–678 (2021). https://doi.org/10.1007/s00382-020-05498-7
873

874   Coxon, G., Addor, N., Bloomfield, J. P., Freer, J., Fry, M., Hannaford, J., Howden, N. J. K., Lane, R.,
875   Lewis, M., Robinson, E. L., Wagener, T., and Woods, R. (2020). CAMELS-GB: hydrometeorological
876   time series and landscape attributes for 671 catchments in Great Britain, Earth Syst. Sci. Data, 12, 2459–
877   2483, https://doi.org/10.5194/essd-12-2459-2020.
878
879   Coppola, E., Nogherotto, R., Ciarlò, J. M., Giorgi, F., van Meijgaard, E., Kadygrov, N., et al.
880   (2021). Assessment of the European Climate Projections as Simulated by the Large EURO-CORDEX
881   Regional and Global Climate Model Ensemble. Journal of Geophysical Research: Atmospheres, 126,
882   e2019JD032356. https://doi.org/10.1029/2019JD032356
883
884   Demargne, J. et al. (2014). The Science of NOAA's Operational Hydrologic Ensemble Forecast
885   Service. Bull. Amer. Meteor. Soc., 95, 79–98, https://doi.org/10.1175/BAMS-D-12-00081.1.
886
887   Feng, D., Fang, K., & Shen, C. (2020). Enhancing streamflow forecast and extracting insights using long-
888   short term memory networks with data integration at continental scales. Water Resources Research, 56,
889   e2019WR026793. https://doi.org/ 10.1029/2019WR026793
890
891   Feng, D., Liu, J., Lawson, K., & Shen, C. (2022). Differentiable, learnable, regionalized process-based
892   models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy. Water
893   Resources Research, 58, e2022WR032404. https://doi.org/10.1029/2022WR032404
894
895   Feng, D., Beck, H., Lawson, K., and Shen, C. (2023). The suitability of differentiable, physics-informed
896   machine learning hydrologic models for ungauged regions and climate change impact assessment,
897   Hydrol. Earth Syst. Sci., 27, 2357–2373, https://doi.org/10.5194/hess-27-2357-2023.
898
899   Frame, J.M., Kratzert, F., Gupta, H.V., Ullrich, P., & Nearing, G.S. (2022). On Strictly enforced mass
900   conservation constraints for modeling the Rainfall-Runoff process. Hydrological Processes, 37, e14847,
901   https://doi.org/10.1002/hyp.14847.
902
903   Frame, J.M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., et al. (2021b). Deep learning
904   rainfall-runoff predictions of extreme events. *Hydrology and Earth System Sciences*, 26, 3377-
905   3392, https://doi.org/10.5194/hess-26-3377-2022.
906
907   Frame, J.M., Kratzert, F., Raney II, A., Rahman, M., Salas, F.R., & Nearing, G.S. (2021a). Post-
908   processing the National Water Model with Long Short-Term Memory networks for streamflow
909   predictions and diagnostics. *Journal of the American Water Resources Association*, 1-12.
910   https://doi.org/10.1111/1752-1688.12964
911
912   Fry, L. M., Hunter, T. S., Phanikumar, M. S., Fortin, V., and Gronewold, A. D. (2013), Identifying
913   streamgage networks for maximizing the effectiveness of regional water balance modeling, Water Resour.
914   Res., 49, 2689– 2700, doi:10.1002/wrcr.20233.
915
916   Gasset, N., Fortin, V., Dimitrijevic, M., Carrera, M., Bilodeau, B., Muncaster, R., Gaborit, É., Roy, G.,
917   Pentcheva, N., Bulat, M., Wang, X., Pavlovic, R., Lespinas, F., Khedhaouiria, D., and Mai, J.: A 10 km
918   North American precipitation and land-surface reanalysis based on the GEM atmospheric model, Hydrol.
919   Earth Syst. Sci., 25, 4917–4945, https://doi.org/10.5194/hess-25-4917-2021, 2021.
920
921   Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., & Hochreiter, S. (2021a). Rainfall-runoff
922   prediction at multiple timescales with a single Long Short-Term Memory network. *Hydrology and Earth
923   System Sciences*, *25*, 2045-2062. https://doi.org/10.5194/hess-25-2045-2021
924

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

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

932   Gordon, B.L., Brooks, P.D., Krogh, S.A., Boisrame, G.F.S., Carrol, R.W.H., McNamara, J.P., & Harpold,
933   A.A. (2022), Why does snowmelt-driven streamflow response to warming vary? A data-driven review
934   and predictive framework, *Environmental Research Letters*, 15 (5), 053004. https://doi.org/10.1088/1748-
935   9326/ac64b4

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

940   Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F. (2009). Decom- position of the mean squared
941   error and NSE performance criteria: Implications for improving hydrological modelling, J. Hydrol., 377,
942   80–91.

944   Hamon, W. R. (1963). Estimating Potential Evapotranspiration, T. Am. Soc. Civ. Eng., 128, 324–
945   338, https://doi.org/10.1061/TACEAT.0008673.

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

951   Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735-
952   1780. https://doi.org/10.1162/neco.1997.9.8.1735

954   Hoedt, P.J., F. Kratzert, D. Klotz, C. Halmich, M. Holzleitner, G. Nearing, et al. (2021). MC-LSTM:
955   Mass-Conserving LSTM. *arXiv e-prints*, arXiv:2101.05186. Retrieved from
956   https://arxiv.org/abs/2101.05186

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

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

965   Kapnick, S., & Hall, A. (2010). Observed Climate–Snowpack Relationships in California and their
966   Implications for the Future, *Journal of Climate*, 23(13), 3446-
967   3456. https://doi.org/10.1175/2010JCLI2903.1

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

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

976
977     Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *arXiv e-prints*,
978     arXiv:1412.6980. Retrieved from https://arxiv.org/abs/1412.6980

979
980     Konapala, G., Kao, S. C., Painter, S., & Lu, D. (2020). Machine learning assisted hybrid models can
981     improve streamflow simulation in diverse catchments across the conterminous US. Environmental
982     Research Letters, 15(10), 104022. https://doi.org/10.1088/1748-9326/aba927
983
984     Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019a).
985     Toward improved predictions in ungauged basins: Exploiting the power of machine learning. *Water*
986     *Resources Research*, *55*, 11,344–11,354. https://doi.org/10.1029/2019WR026065
987
988     Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. S. (2019b). Towards
989     learning universal, regional, and local hydrological behaviors via machine learning applied to large-
990     sample datasets. *Hydrology and Earth System Sciences*, *23*, 5089-5110. https://doi.org/10.5194/hess-23-
991     5089-2019
992
993     Kratzert, F., Klotz, D., Hochreiter, S., & Nearing, G. S. (2021). A note on leveraging in multiple
994     meteorological data sets with deep learning for rainfall-runoff modeling. Hydrology and Earth System
995     Sciences, 25(5), 2685–2703. https://doi.org/10.5194/hess-25-2685-2021.
996
997     Lai, C., Chen, X., Zhong, R., and Wang, Z. (2022), Implication of climate variable selections on the
998     uncertainty of reference crop evapotranspiration projections propagated from climate variables
999     projections under climate change, Agricultural Water Management, 259(1), 107273,
1000    https://doi.org/10.1016/j.agwat.2021.107273.
1001
1002    Lee, D., Lee, G., Kim, S., & Jung, S. (2020). Future Runoff Analysis in the Mekong River Basin under a
1003    Climate Change Scenario Using Deep Learning. *Water*, *12*(6):1556. https://doi.org/10.3390/w12061556
1004
1005    Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., et al. (2021). Hydrological concept
1006    formation inside long short-term memory (LSTM) networks. Hydrology and Earth System Sciences, 26
1007    (12), https://doi.org/10.5194/hess-26-3079-2022.
1008
1009    Lehner, F., Wahl, E., R., Wood, A. W., Blatchford, D. B., & Llewellyn, D. (2017). Assessing recent
1010    declines in Upper Rio Grande runoff efficiency from a paleoclimate perspective. *Geophysical Research*
1011    *Letters*, *44*, 4124-4133. https://doi.org/10.1002/2017GL073253
1012
1013    Lehner, B., Verdin, K., and Jarvis, A. (2008). New Global Hydrography Derived From Spaceborne
1014    Elevation Data, Eos T. Am. Geophys. Un., 89, 93–94.
1015
1016    Lemaitre-Basset, T., Oudin, L., Thirel, G., and Collet, L.: Unraveling the contribution of potential
1017    evaporation formulation to uncertainty under climate change, Hydrol. Earth Syst. Sci., 26, 2147–2159,
1018    https://doi.org/10.5194/hess-26-2147-2022, 2022.
1019
1020    Li, K., Huang, G., Wang, S., Razavi, S., & Zhang, X. (2022). Development of a joint probabilistic
1021    rainfall-runoff model for high-to-extreme flow projections under changing climatic conditions. Water
1022    Resources Research, 58, e2021WR031557. https://doi. org/10.1029/2021WR031557
1023

1024 Lin, L., Gettelman, A., Fu, Q. et al. Simulated differences in 21st century aridity due to different scenarios
1025 of greenhouse gases and aerosols. Climatic Change 146, 407–422 (2018). https://doi.org/10.1007/s10584-
1026 016-1615-3

1027
1028 Liu, X., Li, C., Zhao, T., and Han, L. (2020) Future changes of global potential evapotranspiration
1029 simulated from CMIP5 to CMIP6 models, Atmospheric and Oceanic Science Letters, 13:6, 568-
1030 575, DOI: 10.1080/16742834.2020.1824983

1031
1032 Liu, Z., Han, J., and Yang, H. (2022), Assessing the ability of potential evaporation models to capture the
1033 sensitivity to temperature, Agricultural and Forest Meteorology, 317, 108886.

1034
1035 Liu, Z., Wang T., Han, J., Yang, W., & Yang, H. (2022). Decreases in mean annual streamflow and
1036 interannual streamflow variability across snow-affected catchments under a warming climate.
1037 *Geophysical Research Letters*, *49*(3), e2021GL097442. https://doi.org/10.1029/2021GL097442

1038
1039 Lofgren, B.M., Hunter, T.S., Wilbarger, J. (2011), Effects of using air temperature as a proxy for potential
1040 evapotranspiration in climate change scenarios of Great Lakes basin hydrology, Journal of Great Lakes
1041 Research, 37 (4), 744-752.

1042
1043 Lofgren, B. M., and Rouhana, J. (2016) Physically Plausible Methods for Projecting Changes in Great
1044 Lakes Water Levels under Climate Change Scenarios. J. Hydrometeor., 17, 2209–
1045 2223, https://doi.org/10.1175/JHM-D-15-0220.1.

1046
1047 Lu, D., Konapala, G., Painter, S. L., Kao, S. C., & Gangrade, S. (2021). Streamflow simulation in data-
1048 scarce basins using Bayesian and physics-informed machine learning models. Journal of
1049 Hydrometeorology, 22(6), 1421– 1438. https://doi.org/10.1175/JHM-D-20-0082.1

1050
1051 Lu, J., Sun, G., McNulty, S.G. and Amatya, D.M. (2005), A comparison of six potential
1052 evapotranspiration methods for regional use in the southeastern United States. JAWRA Journal of the
1053 American Water Resources Association, 41: 621-633. https://doi.org/10.1111/j.1752-
1054 1688.2005.tb03759.x

1055
1056 Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., et al. (2021). Transferring hydrologic
1057 data across continents – leveraging data-rich regions to improve hydrologic prediction in data-sparse
1058 regions. Water Resources Research, 57, e2020WR028600. https://doi. org/10.1029/2020WR028600

1059
1060 Mai et al. (2022). The Great Lakes runoff intercomparison project phase 4: the Great Lakes (GRIP-GL),
1061 Hydrologic and Earth System Sciences, 26 (13), 3537-3573, https://doi.org/10.5194/hess-26-3537-2022.

1062
1063 Martens, B., Miralles, D. G., Lievens, H., van der Schalie, R., de Jeu, R. A. M., Fernández-Prieto, D.,
1064 Beck, H. E., Dorigo, W. A., and Verhoest, N. E. C. (2017). GLEAM v3: satellite-based land evaporation
1065 and root-zone soil moisture, Geosci. Model Dev., 10, 1903– 1925, https://doi.org/10.5194/gmd-10-1903-
1066 2017.

1067
1068 Martin, J. T., Pederson, G. T., Woodhouse, C. A., Cook, E. R., McCabe, G. J., Anchukaitis, K. J., et al.
1069 (2020). Increased drought severity tracks warming in the United States' largest river basin. *Proceedings
1070 of the National Academy of Sciences*, *117*(21). https://doi.org/10.1073/pnas.1916208117

1071
1072 McCabe, G. J., Wolock, D. M., Pederson, G. T., Woodhouse, C. A., & McAfee, S. (2017). Evidence that
1073 recent warming is reducing upper Colorado River flows. *Earth Interactions*, *21*(10), 1-14.
1074 https://doi.org/10.1175/EI-D-17-0007.1

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

Milly, P. C. D., & Dunne, K. A. (2020). Colorado River flow dwindles as warming-driven loss of reflective snow energizes evaporation. *Science*, *367*(6483), 1252-1255. https://doi.org/10.1126/science.aay9187

Monteith, J. L. (1965), Evaporation and environment, in: Symposia of the society for experimental biology, volume 19, Cambridge University Press (CUP), Cambridge, UK, 205–234 pp.

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

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

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

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

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

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

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

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

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

Rungee, J., Ma, Q., Goulden, M. L., & Bales, R. (2021). Evapotranspiration and runoff patterns across California's Sierra Nevada. *Frontiers in Water*, *3:655485*. https://doi.org/10.3389/frwa.2021.655485

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

1125
1126 Shaw, S.B. and Riha, S.J. (2011), Assessing temperature-based PET equations under a changing climate
1127 in temperate, deciduous forests. Hydrol. Process., 25: 1466-1478. https://doi.org/10.1002/hyp.7913
1128
1129 Steinman, A.D. et al. (2017), Ecosystem services in the Great Lakes, Journal of Great Lakes Research, 43
1130 (3), 161-168. https://doi.org/10.1016/j.jglr.2017.02.004
1131
1132 Stewart, I. T., Cayan, D. R., & Dettinger, M. D. (2005). Changes toward Earlier Streamflow Timing
1133 across Western North America, *Journal of Climate*, 18(8), 1136-1155.
1134 https://doi.org/10.1175/JCLI3321.1
1135
1136 Szilagyi, J., Crago, R., and Qualls, R. (2017), A calibration-free formulation of the complementary
1137 relationship of evaporation for continental-scale hydrology, J. Geophys. Res. Atmos., 122, 264– 278,
1138 doi:10.1002/2016JD025611.
1139
1140 Taranu, I.S., Somot, S., Alias, A. et al. Mechanisms behind large-scale inconsistencies between regional
1141 and global climate model-based projections over Europe. Clim Dyn 60, 3813–3838 (2023).
1142 https://doi.org/10.1007/s00382-022-06540-6
1143
1144 Wang, Q. J. (1991). The genetic algorithm and its application to calibrating conceptual rainfall-runoff
1145 models, *Water Resources Research*, *27*(9), 2467-2471. https://doi.org/10.1029/91WR01305
1146
1147 Wi, S., & Steinschneider, S. (2022). Assessing the physical realism of deep learning hydrologic model
1148 projections under climate change. Water Resources Research, 58,
1149 e2022WR032123. https://doi.org/10.1029/2022WR032123
1150
1151 Woodhouse, C. A., & Pederson, G. T. (2018). Investigating runoff efficiency in upper Colorado river
1152 streamflow over past centuries. *Water Resources Research*, *54*, 286-300.
1153 https://doi.org/10.1002/2017WR021663
1154
1155 Yilmaz, K. K., Gupta, H. V., and Wagener, T. (2008). A process-based diagnostic approach to model
1156 evaluation: Application to the NWS distributed hydrologic model, Water Resour. Res., 44, 1–18.
1157
1158 Zhong, L., Lei, H., & Gao, B. (2023). Developing a physics-informed deep learning model to simulate
1159 runoff response to climate change in Alpine catchments. Water Resources Research, 59,
1160 e2022WR034118. https://doi.org/10.1029/2022WR034118
1161