

Response to Reviewers for: ‘On the need for physical constraints in deep learning rainfall-runoff projections under climate change: a sensitivity analysis to warming and shifts in potential evapotranspiration’

Sungwook Wi, Scott Steinschneider

Corresponding author: Sungwook Wi, sw2275@cornell.edu

Key

Black font: Reviewer comments

Blue font: Author responses

Italicized orange font: Updated manuscript wording, underline for changes to original

We greatly appreciate the time and detail that the three reviewers put into their evaluation of our manuscript. We have addressed all of the comments on a point-by-point basis, which we detail below. In particular, we have made substantial revisions throughout the manuscript to address several major concerns expressed by the reviewers, including:

- A more nuanced presentation of the specific experiments conducted in this work, emphasizing that we conduct sensitivity analyses rather than formal projections under climate change;
- Improved context for our results and experiments, in particular with respect to the machine learning literature on distribution shifts and causality;
- More detail on how our results highlight the strength and weaknesses of DL models to simulate hydrologic responses under changing snow accumulation and melt processes with warming; and
- A more thorough treatment of parametric uncertainty in the process models, and how it impacts the interpretation of DL model responses to warming.

We think these revisions, along with several others, have served to significantly improve the manuscript.

Reviewer #1

This study examines the behavior of different models under a hypothetical scenario where 4°C are added to the daily minimum and maximum temperatures. In doing so, the contribution finds that models with more explicit representations of hydrological processes are likely to exhibit more realistic behaviors under this shift.

I find this kind of study very important and very timely (as other discussion papers show; see e.g.: Reichert et al. 2023). On top of that, the execution is done well: The work is, by and large, well motivated; the idea is good; and, all tables and images are clear; almost everything is documented. I therefore think that the study should definitely be published on HESS. In terms of critique I have one point about the literature that I think is crucial, and some small questions/comments. The latter are, however, not so important.

We thank the reviewer for their overall positive, constructive, and speedy review. We greatly appreciate the feedback and believe it has served to significantly improve the manuscript.

Major Comment

The references are quite thorough with regard to the recent use of deep learning in hydrology. I complement the authors for that. They do, however, ignore large amounts of work from the outside the field. Normally this would not be a concern --- since one feeds into the other --- but here it does skew the motivation somewhat. As of now the introduction/motivation of the work reads as if current researcher are not aware that one can increase the temperature by some degrees and then test what the model would do under such circumstances. This is however not the case. For example, the group I am involved with, did not conduct such counterfactual experiments because we knew that deep learning models are out of the box not be able cope with arbitrary shifts in the covariance structures of the inputs. Statistical learning hinges on the idea that the future looks similar to the past --- and in a counterfactual setting this property is not given by design.

I strongly believe that the paper should give a better overview of the current machine learning literature and use that to discuss the merits and limits of the study design. This would give readers a much richer picture of what the proposed evaluation can probe.

Specifically, I am thinking that the paper should reference current work on (a) causality and (b) distribution shifts; and then use it feed into the discussion of the limitations of the current study. The reason why I think of (a) and (b) is that both research branches are fundamental to understand the study design: (a) Causality is important because the examination is a true counterfactual in that the adopted input has not --- and will never be --- observed in reality (remember, the daily values of the min and max temperatures change by adding exactly 4°C to all basins, while inputs like the radiation, wind, precipitation, and vapor pressure remain entirely the same). (a) The research on distribution shifts is important because adding 4°C to each day is a prime example of a covariate shift. Detecting, handling, "robustifying" and/or adapting to distribution shifts is an active area of research and should be seen as an open problem. Roughly speaking, results from (a) and (b) provide a counter point to the current motivation of the research in that they suggest that data-driven models should per-se not be able to withstand a counterfactual examination. I think this would help readers to understand that the "physical plausible" response of the catchment model is measured with a "physically implausible" counterfactual signal (which is not observed in any catchment no matter what and will force the models into a sort of "extrapolation regime"). I believe that only then readers will understand that this is a very special form of test --- and that is very impressive that it is possible to design data-driven models that already show promising result in this setting, while having just a few more inductive biases than the current LSTM based rainfall-runoff models. In this regard, I do not want to force the authors to cite any particular work, but beg them to align their work with these branches of research (even if it means that they need to relativize their a-priori expectations)

We are grateful to the reviewer for making this suggestion. We agree that the literature on causality and distribution shifts in machine learning is extremely relevant to our study design, and in particular to its interpretation, limitations, as well as a fertile ground for future work. The

last two concluding paragraphs in our original manuscript were an attempt to address some of these issues, although admittedly this was not done to the extent necessary or with reference to the large body of work on these topics in the broader ML literature that the reviewer notes here. Consistent with one of the last reviewer suggestions below, we have taken this opportunity to significantly revise (i.e., large rewrite) our Discussion and Conclusion section, removing some of the older content and replacing it with a more robust treatment of the issues raised here. In the process of this revision, we tried to integrate our discussion around physics-informed machine learning (PIML) into the broader discussion on distribution shifts and causality, as we view PIML as one set of approaches (among several) that falls under the broader umbrella of causal deep learning methods.

We also note that in response to this comment and another by Reviewer #2, we have significantly revised the text throughout our manuscript to better convey what our experiment actually tests: the sensitivity of these models to imposed shifts in temperature and associated changes in potential evapotranspiration, rather than internally consistent climate changes across all meteorological variables. We believe these revisions also support the general points being made in this reviewer's comment.

Minor Comments

L. 85-86. Please add a reference to this sentence (or an explanation why no reference is given). You make the claim that "many argue" without even giving a single example.

We have added a recent paper (Nearing et al., 2022) that makes this argument based on past literature (or an adjusted version of this argument, see response to comment below), and also changed "many" to "some" to avoid overstating this claim.

L. 85-86. I think the meaning of "state-of-the-science" should be outlined. As far as I am aware it is not common terminology in hydrology (I, for one, had to look it up and am still not sure what is meant with it in this context).

We have changed the wording here to be more explicit, removing 'state-of-the-science' and instead replacing it with 'most accurate and extrapolatable'

L.100-101. I disagree with the claim about the corollary. Maybe it is an implication? I am not sure however: (a) Given the noise in the data, even without new climate conditions the predictions might be physically implausible. (b) Just because a ML model is "physically plausible" in in a out os sample setting does not mean that it remains so under a shift setting. What do you think about writing something like "From these results one might think that ..." or "If we spin these results further one could think that...".

We agree a wording change is warranted here. From the suggestions provided, we have altered this sentence to read:

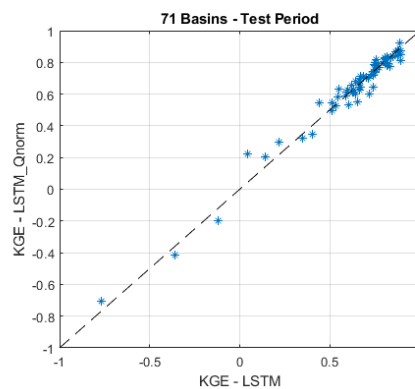
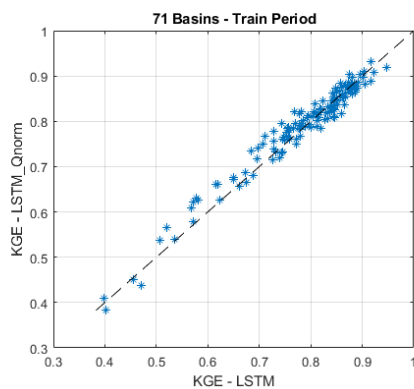
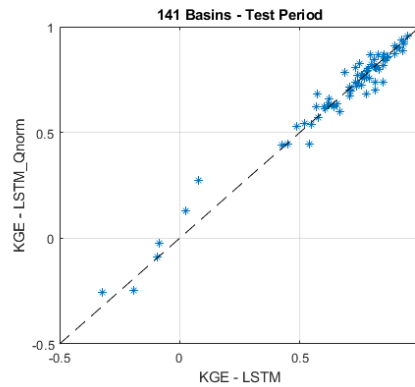
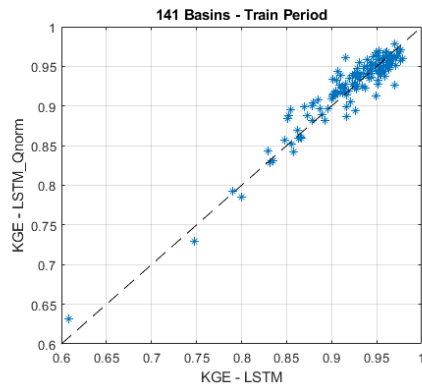
A potential implication of this finding might be that these models can produce physically plausible streamflow predictions under new climate conditions.

L.108ff Is it correct that, from a hydrological perspective, this assume that there are no Glaciers/permanent snow in the basins (which, I think is, e.g., not true for CAMELS US as used in Liu et al. 2022)? The mechanism would be that as long as there is more melting happening, we should see higher water levels with higher temperatures.

That is correct, and our study in Wi and Steinschneider (2022) does point out (and highlight in the results) that this assumption of streamflow loss under warming would not apply to watersheds that drain regions with glaciers or over-year snow cover. We have adjusted this line to highlight this exception.

L. 334. Would you be so kind and make a comparison of the (normal) LSTM performance with normalized streamflow and without it? I once made a similar test, where I trained an LSTM on CMALES US without setting the standard deviations to one. And, got very bad results... However, your performance seems to be comparable to the ones reported in Mai et al. (2022). To me it is not really clear how you did that (especially since you used a relatively small learning rate and since the linear layer requires much bigger parameters in your setting). Maybe it is because the magnitude and behavior of the GRIP-GL rivers are much less diverse than the ones in CAMELS US?

Sure thing. We went ahead and refit the LSTM using normalized streamflow, employing the exact same training process as for the LSTM described in the article. We've pasted below a comparison of performance (KGE) between the original LSTM (without normalization) and the LSTM with normalization, showing results separately for training and testing sites and training and testing periods. We do not see any meaningful difference between the two. We agree that perhaps this result is driven by the fact that the diversity across sites in our domain is less than across the entire continuous US. Related to this, the flow in our region is much less skewed compared to some other arid and semi-arid regions, which might be contributing to this result.



L. 165ff. I know this is a choice of style and I will not mention this for the other occurrences, but: I would appreciate if you could already sketch the outcome of the experiments here (and in the other instances where you hypothesize about properties that one actually already knows at the time of writing).

We have adjusted the wording throughout the Introduction to state the outcomes we found, rather than posit them as hypothesized outcomes.

L. 176. Maybe adjust sentence a bit. I pretty sure that Frame et al. 2022 did not made an argument that physical constraints are not needed in for generating plausible projections under climate change. And, this sentence could easily be misread in that way.

To avoid any confusion or misinterpretation, we have simply removed the reference to Frame et al. 2022 (and arguments that physics-informed constraints are unneeded) altogether. We think the sentence stands fine on its own without this added clause.

L.268ff & L.350-351. It is probably an oversight on my side, but cannot find the code for this analysis in the zenodo repository.

This was an oversight on our part. The code for the National LSTM has now been added to the Zenodo repository.

L.344ff. Can you add a description or table with the grid you searched the hyper-parameters for to the supplementary?

Yes, we have added a table in the Supporting Information that shows the grid search used for the hyper-parameters, and now reference this SI table in the main manuscript.

L.377. I would recommend to explicitly write about σ and $\hat{\sigma}$ here so that readers know what you are referring to.

We now include the equation for $\hat{\sigma}(\cdot)$ in relation to $\sigma(\cdot)$, and provide a brief explanation of its output.

Table 2. I think the MC-LSTM KGE for "Testing Sites: Testing Period" should also be marked in bold since it is also 0.72 (the decimals that follow and are not shown should not be considered for a tie breaker here).

We agree, we have bolded the 0.72 being referred to in this table. We also slightly revised the manuscript text to be consistent with this change.

In particular, the LSTM outperforms the MC-LSTM and MC-LSTM-PET for NSE and FLV (as well as KGE in the training period), the MC-LSTM-PET outperforms the LSTM and MC-LSTM for PBIAS, and either the MC-LSTM or MC-LSTM-PET are the best performers for FHV.

L.497ff Please describe the actual changes that you made to the static attributes either here or in the supplementary. I can see the changes in the data, but that requires readers to reconstruct what you did.

We have added a section to our Supporting Information to more clearly describe the changes made to the static attributes, and refer to this section here in the main article.

L.497ff I am probably missing something here, but to me its is not obvious why you changed the snow fraction of the precipitation with temperatures below 0°C? If the model gets an input with -3°C it should not matter to this whether this value was the true input or the counterfactually modified one; no?

The static input frac_snow is defined as the fraction of precipitation falling on days with mean daily temperatures below 0°C, i.e., the total amount of precipitation falling on days with $T < 0C$ divided by the total amount of precipitation falling on all days. Under our warming scenario, the number of days with precipitation falling when temperatures are below 0°C declines, and thus, so does frac_snow. We now clarify this in our revised manuscript (see response to comment directly above).

L.656 consist -> consistent

This has been corrected.

L. 803ff. Is it really necessary to discuss short-wave radiation for so long here? You also did not consider that the thermic and dynamic behavior of the atmosphere and hence, the precipitation patterns would, for example, change over the whole region. I think you could abbreviate this paragraph considerably by just stating that the input modification is pragmatic and intuitiv, but does not reflect how the meteorological behavior would actually play out under climate change. This would then also my proposed literature references if you decide to include it.

We have taken the reviewer's suggestion, and have significantly shortened our focus on radiation here in favor of a broader treatment of the issues of distribution shifts and causality, as mentioned in the reviewer's main comment.

Upon reflection I would like to add that I think it would be highly beneficial if you could add some representative Hydrographs to an Appendix. This is, for one because I am interested to see some because of my personal experience with mass-conserving models; but secondly I also genuinely believe that it would help readers to put the performance and interventions into perspective.

When describing the results in Figure 7, we now reference individual hydrographs for specific sites (at both daily and monthly timescales), which are provided in the SI. We reference these SI figures while highlighting changes to key attributes of streamflow (FLV, FHV, COM) under warming, in an effort to better show what some of these differences in flow statistics mean in terms of daily flow time series.

References

- Reichert, P., Ma, K., Höge, M., Fenicia, F., Baity-Jesi, M., Feng, D., and Shen, C.: Metamorphic Testing of Machine Learning and Conceptual Hydrologic Models, *Hydrol. Earth Syst. Sci. Discuss.* [preprint], <https://doi.org/10.5194/hess-2023-168>, in review, 2023.
- Mai, J., Shen, H., Tolson, B. A., Gaborit, É., Arsenault, R., Craig, J. R., Fortin, V., Fry, L. M., Gauch, M., Klotz, D., Kratzert, F., O'Brien, N., Princz, D. G., Rasiya Koya, S., Roy, T., Seglenieks, F., Shrestha, N. K., Temgoua, A. G. T., Vionnet, V., and Waddell, J. W.: The Great Lakes Runoff Intercomparison Project Phase 4: the Great Lakes (GRIP-GL), *Hydrol. Earth Syst. Sci.*, 26, 3537–3572, <https://doi.org/10.5194/hess-26-3537-2022>, 2022.

1 **On the need for physical constraints in deep learning rainfall-runoff**
2 **projections under climate change: a sensitivity analysis to warming and shifts**
3 **in potential evapotranspiration**

4
5 **Sungwook Wi¹, Scott Steinschneider¹**

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

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

Abstract

Deep learning (DL) rainfall-runoff models ~~have recently emerged as state-of-the-science tools for hydrologic prediction that~~ outperform ~~conventional~~, process-based models in a range of applications. However, it remains unclear whether ~~deep learning~~DL models can produce physically plausible projections of streamflow under ~~significant amounts of~~ climate change. We investigate this question ~~here, focusing specifically on~~through a sensitivity analysis of modeled responses to increases in temperature and potential evapotranspiration (PET), ~~with other meteorological variables left unchanged~~. Previous research has shown that temperature-based ~~PET~~ methods ~~to estimate PET lead to~~ overestimates ~~evaporative of~~ water loss ~~in rainfall-runoff models~~ under warming, as compared to energy budget-based PET methods. ~~We therefore~~ Consequently, we assess the reliability of streamflow projections under warming by comparing projections ~~with both temperature-based and energy budget-based PET estimates,~~ assuming that reliable streamflow ~~projections responses to warming~~ should exhibit less ~~evaporative~~ water loss when forced with smaller, ~~(energy budget-based)~~ ~~projections of future~~ PET ~~compared to temperature-based PET~~. We conduct this assessment using three ~~conceptual~~, process-based ~~rainfall-runoff models~~ ~~rainfall-runoff models~~ and three ~~deep learning~~DL models, trained and tested across 212 watersheds in the Great Lakes basin. The ~~deep learning~~DL models include a ~~regional~~ Long Short-Term Memory network (LSTM), a mass-conserving LSTM (MC-LSTM) ~~that preserves the water balance~~, and a novel variant of the MC-LSTM that also respects the relationship between PET and ~~evaporative~~ water loss (MC-LSTM-PET). ~~After validating models against~~ We first ~~compare~~ historical streamflow ~~and actual watershed-scale evapotranspiration, predictions from all models under spatial and temporal validation, and also assess model skill in estimating watershed-scale evapotranspiration.~~ We then ~~we~~ force all models with scenarios of warming, historical precipitation, and both temperature-based (Hamon) and energy budget-based (Priestley-Taylor) PET, and compare their ~~projections responses for changes in average long-term mean daily~~ flow, as well as low flows, high flows, and ~~seasonal~~ streamflow ~~seasonal~~ timing. ~~Finally, w~~We also explore similar ~~projections~~

51 responses using a National LSTM fit ~~to to a broader set of~~ 531 watersheds across the ~~contiguous~~ United
52 States to assess how the inclusion of a larger and more diverse set of basins influences signals of hydrologic
53 response under warming. The main results of this study are as follows:

- 54 1. The three Great Lakes ~~deep learning~~DL models ~~significantly~~ substantially outperform all process
55 models in streamflow estimation ~~under spatiotemporal validation, with only small differences~~
56 ~~between the DL models~~. The MC-LSTM-PET also matches the best process models and
57 outperforms the MC-LSTM in estimating actual evapotranspiration ~~under spatiotemporal~~
58 ~~validation~~.
- 59 2. All process models show a downward shift in long-term mean daily average flows under warming,
60 but ~~this median shifts~~ is ~~is~~ are ~~significantly~~ considerably larger under temperature-based PET (17%
61 to 25%) ~~estimates~~ than energy budget-based PET (-6% to -9%). The MC-LSTM-PET model
62 exhibits similar differences in water loss across the different PET forcings, ~~consistent with the~~
63 ~~process models~~. ~~However~~ Conversely, the LSTM exhibits unrealistically large water losses under
64 warming ~~as compared to the process models~~ using Priestley-Taylor PET (20%), while the MC-
65 LSTM is relatively insensitive to PET method.
- 66 3. ~~All deep learning~~DL models exhibit smaller changes in high flows and streamflow seasonal timing
67 of flows as compared to the process models while ~~deep learning~~DL ~~projections estimates~~
68 flows are ~~all very consistent and~~ within the range projected estimated by the process models.
- 69 4. Like the Great Lakes LSTM, the National LSTM also shows unrealistically large water losses under
70 warming (25%), but ~~.~~ ~~However, when compared to the Great Lakes deep learning models,~~
71 ~~projections from the National LSTM were~~ it is more stable when many inputs ~~were~~ are changed
72 under warming and better aligned ~~sed~~ with process model ~~projections responses~~ for streamflow
73 seasonal timing of flows. ~~This suggests that the addition of more, diverse watersheds in training~~
74 ~~does help improve climate change projections from deep learning models, but this strategy alone~~
75 ~~may not guarantee reliable projections under unprecedented climate change.~~

76 Ultimately, the results of this [work-sensitivity analysis](#) suggest that physical considerations regarding model
77 architecture and input variables ~~are~~ may be necessary to promote the physical realism of deep learning-
78 based hydrologic projections under climate change.

79

80 **Keywords**

81 Deep learning, machine learning, Long Short-Term Memory network, LSTM, Great Lakes, climate
82 change, rainfall-runoff

83

84

85

86

87

88

89

90

91

92 **1. Introduction**

93 Rainfall-runoff models are used throughout hydrology in a range of applications, including retrospective
94 streamflow estimation (Hansen et al. 2019), streamflow forecasting (Demargne et al., 2014), and prediction
95 in ungauged basins (Hrachowitz et al., 2013). Work over the last few years has demonstrated that deep
96 learning (DL) rainfall-runoff models (e.g., Long Short-Term Memory networks (LSTMs); Hochreiter and
97 Schmidhuber, 1997) outperform conventional process-based models in each of these applications,
98 especially when those DL models are trained with large datasets collected across watersheds with diverse
99 climates and landscapes (Kratzert et al., 2019a,b; Feng et al., 2020; Ma et al., 2021; Gauch et al., 2021a,b;
100 Nearing et al., 2021). For example, in one extensive benchmarking study, Mai et al. (2022) found that a

101 regionally trained LSTM outperformed 12 other lumped and distributed process-based models of varying
102 complexity in rivers and streams throughout the Great Lakes basin. These and similar results have led ~~many~~
103 some to argue that DL models represent the most accurate and extrapolatable rainfall-runoff models
104 available (Nearing et al., 2022)~~most accurate and extrapolatable rainfall-runoff models available (Nearing~~
105 ~~et al., 2022)~~.

106
107 However, there remains one use case of rainfall-runoff models where the superiority of DL is unclear: long-
108 term projections of streamflow under climate change. Past studies using DL rainfall-runoff models for
109 hydrologic projections under climate change are rare (Lee et al., 2020; Li et al., 2022), and few have
110 evaluated their physical plausibility (Razavi, 2021; Reichert et al., 2023; Zhong et al., 2023). A reasonable
111 concern is whether DL rainfall-runoff models can extrapolate hydrologic response under unprecedented
112 climate conditions, given that they are entirely data driven and do not explicitly represent the physics of the
113 system. It is not clear *a priori* whether this concern has merit, because DL models fit to a large and diverse
114 set of basins have the benefit of learning hydrologic response across climate and landscape gradients. In so
115 doing, the model can, for example, learn hydrologic responses to climate in warmer regions and then
116 transfer this knowledge to projections of streamflow in cooler regions subject to climate change induced
117 warming. In addition, past work has shown that LSTMs trained only to predict streamflow have memory
118 cells that strongly correlate with independent measures of soil moisture and snowpack (Lees et al. 2022⁺),
119 suggesting that DL hydrologic models can learn fundamental hydrologic processes. A ~~corollary potential~~
120 implication ~~to of~~ this finding ~~is might be~~ that these models can ~~may~~ produce physically plausible streamflow
121 predictions under new climate conditions.

122
123
124
125 It is challenging to assess the physical plausibility of DL-based hydrologic projections under significantly
126 substantially different climate conditions, because there are no future observations against which to

127 compare. This challenge is exacerbated by significant uncertainty in process model projections under
128 alternative climates, which makes establishing reliable benchmarks difficult. Future process model-based
129 projections can vary widely due to both parametric and structural uncertainty (Bastola et al., 2011; Clark et
130 al., 2016; Melsen et al., 2018), and even for models that exhibit similar performance under historical
131 conditions (Krysanova et al., 2018). Assumptions around stationary model parameters are not always valid
132 (Merz et al., 2011; Wallner and Haberlandt, 2015), and added complexity for improved process
133 representation is not always well supported by data (Clark et al., 2017; Towler et al., 2023; Yan et al., 2023).
134 Together, these challenges highlight the difficulty in establishing good benchmarks of hydrologic response
135 under alternative climates against which to compare and evaluate DL-based hydrologic projections under
136 climate change.

137
138
139 Recently, Wi and Steinschneider (2022) (hereafter WS22) ~~addressed this challenge directly,~~
140 ~~forwarding~~forwarded an experimental design to evaluate the physical plausibility of DL hydrologic
141 responses to new climates, in which DL hydrologic models ~~fit to 15 watersheds in California and 531~~
142 ~~catchments across the United States~~ were forced with historical precipitation and temperature, but with
143 temperatures adjusted by up to 4°C. Based on past literature ~~(Cayan et al., 2001; Stewart et al., 2005;~~
144 ~~Kapnick and Hall, 2010; Lehner et al., 2017; McCabe et al., 2017; Dierauer et al., 2018; Mote et al., 2018;~~
145 ~~Woodhouse & Pederson, 2018; Martin et al., 2020; Milly & Dunne, 2020; Rungee et al., 2021; Gordon et~~
146 ~~al., 2022; Liu et al., 2022),~~ WS22 posited that in non-glaciated regions, physically plausible hydrologic
147 ~~projections-responses~~ should show an increase in water loss, defined as water that enters the watershed via
148 precipitation but never contributes to streamflow because it is ‘lost’ to a terminal sink. Specifically, WS22
149 assumed that evaporative water loss should increase and annual ~~decline in total annual~~ average streamflow
150 should decline compared to a baseline ~~historical~~ simulation, due to increases in potential evapotranspiration
151 (PET) with warming (and no changes in precipitation). Results showed that ~~the an~~ LSTM trained to the 15
152 watersheds in California often led to misleading increases in annual runoff under ~~significant~~ warming, while

153 this phenomenon was less likely (though still present) in ~~the a DL~~ model trained to 531 catchments [across](#)
154 [the United States](#).

155
156 WS22 also conducted their experiment with physics-informed machine learning (PIML) models; ~~in which~~
157 ~~data-driven techniques are imbued with process knowledge constructs~~ (Karpatne et al., 2017), ~~WS22~~
158 ~~focused on two PIML strategies for the smaller case study in California~~, using process model output (~~e.g.,~~
159 ~~soil moisture, evapotranspiration (ET)~~) directly as input to the LSTM (similar to Konapala et al., 2020; Lu
160 et al., 2021; Frame et al., 2021a); ~~and also or~~ as additional target variables in a multi-output architecture.
161 The former approach had some success in removing instances of increasing runoff ratio with warming, ~~but~~
162 ~~although~~ this ~~depended heavily on the accuracy of~~ ~~was dependent on the~~ ~~the~~ process ~~model~~ ~~used~~ ~~ET~~.

163
164 Other PIML approaches that more directly adjust the architecture of DL rainfall-runoff models may be
165 better suited for improving long-term streamflow projections under climate change without requiring an
166 accurate process-based model. For instance, Hoedt et al. (2021) introduced a mass conserving LSTM (MC-
167 LSTM) that ensures cumulative streamflow predictions do not exceed precipitation inputs. [Hybrid models](#)
168 [present a related approach, where DL modules are embedded within process models structures \(Jiang et al.,](#)
169 [2020; Feng et al., 2022; Hoge et al., 2022; Feng et al., 2023a\)](#). ~~In some cases, These~~ ~~is~~ ~~architectural~~ ~~changes~~
170 ~~can slightly degrade performance compared to underperformed~~ a standard LSTM ~~when predicting out of~~
171 ~~sample extreme events~~ (Frame et al., 2021b; [Feng et al., 2023b](#)), [but other times such changes can be](#)
172 [beneficial \(Feng et al., 2023a\)](#). ~~and s~~ Some have argued that these physical constraints may inhibit the ability
173 of DL models to learn biases in forcing data (Frame et al. 2022). ~~Still,~~ ~~but~~ the benefits of ~~this such~~ mass
174 conserving architectures ~~s~~ have not been tested when employed under previously unobserved climate change.

175
176 For all models considered in WS22, a major focus was evaluating the direction of annual total runoff change
177 in the presence of warming and no change in precipitation. However, that study did not consider the
178 magnitude of runoff change and how it relates to projected changes in PET. As we argue below, this

179 comparison provides a unique way to assess the physical plausibility of future hydrologic projections.
180 Several studies have investigated the effects of different PET estimation methods on the magnitude of PET
181 and runoff change in a warming climate (Lofgren et al., 2011; Shaw and Riha, 2011; Lofgren and Rouhana,
182 2016; Milly and Dunne, 2017; Lemaitre-Basset et al. 2022). Broadly, ~~this~~ these work studies ~~have~~ shown
183 that temperature-based PET estimation methods (e.g., Hamon, Thornthwaite) ~~significantly~~ substantially
184 overestimate increases in PET under warming as compared to energy budget-based PET estimation methods
185 (e.g., Penman-Monteith, Priestley-Taylor), and consequently lead to unrealistic declines in streamflow
186 under climate change. This is because the actual drying power of the atmosphere is driven by the availability
187 of energy at the surface from net radiation, the current moisture content of the air, temperature (and its
188 effect on the water holding capacity of the air and vapor pressure deficit), and wind speeds. Energy budget-
189 based methods, while imperfect and at times empirical (Greve et al. 2019; Liu et al., 2022), account for
190 some or all of these factors in ways that are generally consistent with their causal impact on PET, while
191 temperature-based methods estimate PET using strictly empirical relationships based largely or entirely on
192 temperature. The latter approach works sufficiently well for rainfall-runoff modeling under historical
193 conditions because of the strong correlation between temperature, net radiation, and PET on seasonal
194 timescales, even though this correlation weakens considerably at shorter timescales (Lofgren et al., 2011).
195 Under climate change, consistent and prominent increases are projected for temperature, but projected
196 changes are less prominent or more uncertain for other factors affecting PET (Lin et al., 2018; Pryor et al.,
197 2020, Liu et al. 2020). Consequently, temperature-based PET methods ~~significantly~~ substantially
198 overestimate future projections of PET compared to energy budget-based methods (Lofgren et al., 2011;
199 Shaw and Riha, 2011; Lofgren and Rouhana, 2016; Milly and Dunne, 2017; Lemaitre-Basset et al. 2022).

200

201 As argued by Lofgren and Rouhana (2016), the bias in PET and runoff that results from different PET
202 estimation methods under warming provides a unique opportunity to assess the physical plausibility of
203 hydrologic projections under climate change. In this study, we adopt this strategy for DL rainfall-runoff
204 models ~~and forward an experimental design~~ through a sensitivity analysis in which both conceptual, process-

205 based and DL hydrologic models are trained with either temperature-based or energy budget-based
206 estimates of PET, along with other meteorological data (precipitation, temperature). These models are then
207 forced with the historical precipitation and temperature series, but with the temperatures warmed by an
208 additive factor and PET calculated from the warmed temperatures using both PET estimation methods. We
209 ~~anticipate~~ show that the process models 1) ~~will~~ exhibit similar performance in historical training and testing
210 periods when using either temperature-based or energy budget-based PET estimates; but 2) ~~will~~ exhibit
211 ~~significantly~~ substantially larger long-term mean streamflow declines under warming when using future
212 PET estimated with a temperature-based method. If the DL rainfall-runoff models follow the same pattern,
213 this would suggest that these models are able to learn the role of PET on evaporative water loss. However,
214 if DL-based models estimate similarly ~~and~~ large long-term mean streamflow declines regardless of the
215 method used to estimate and project PET, this would suggest that the DL models did not learn a mapping
216 between PET and evaporative water loss. Rather, the DL models learned the historical (but non-causal)
217 correlation between temperature and evaporative water loss, and then incorrectly extrapolated that effect
218 into the future with warmer temperatures. ~~Such~~ We show this latter an outcome to be the case, -would ~~which~~
219 indicates that some degree of PIML ~~is~~ may be necessary to guide a DL model towards physically plausible
220 projections under climate change, ~~in contrast to previous arguments against the need for such physical~~
221 constraints (Frame et al. 2022).

222
223 We conduct the experiment above in a case study on 212 watersheds across the Great Lakes basin, using
224 both standard and PIML-based LSTMs. We ~~hypothesize~~ show that a standard LSTM ~~will~~ produces
225 unrealistic hydrologic ~~projections~~ responses to warming because it relies on historical and geographically
226 pervasive correlations between temperature and PET to ~~project~~ estimate streamflow losses ~~under warming~~.
227 We also ~~hypothesize~~ show that PIML-based DL models ~~will~~ bear better able to relate ~~future projections~~
228 of changes in temperature and PET to streamflow change, especially those PIML approaches that directly
229 map PET to evaporative water loss in their architecture.

230

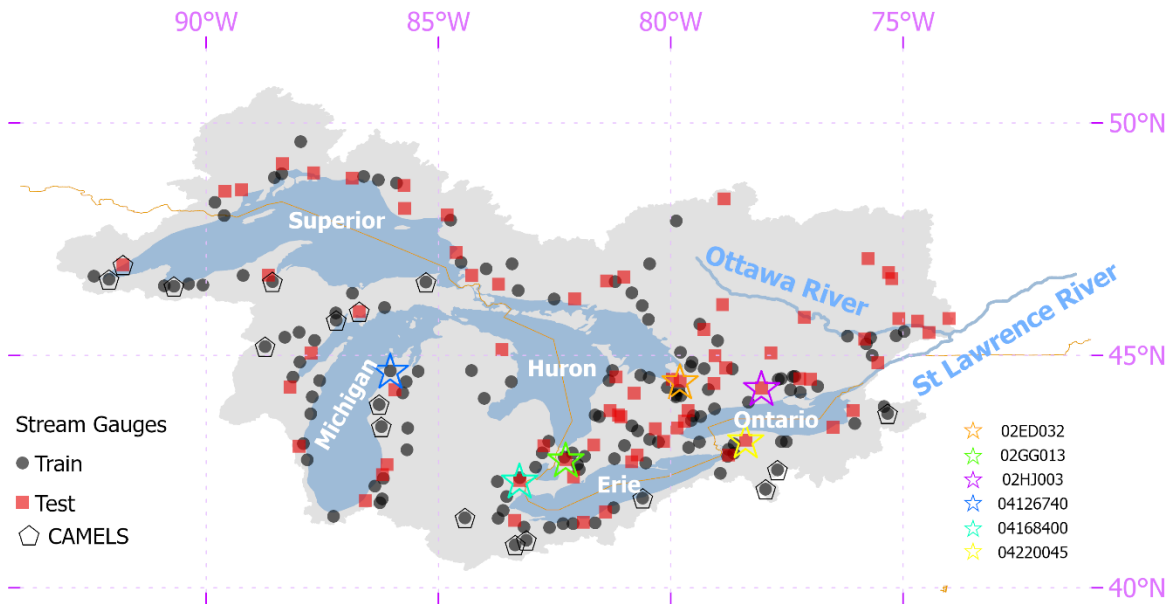
231 The primary goal of this work is to forward an experimental design that can be used to evaluate the
232 suitability of DL rainfall-runoff models for hydrologic projections under climate change, in line with a
233 recent call to design benchmarking studies that assess whether models are fit for specific purposes (Beven,
234 2023). The Great Lakes provides an important case study for this work, given their importance to the culture,
235 ecosystems, and economy of North America (Campbell et al., 2015; Steinman et al., 2017). Projections of
236 future water supplies and water levels in the Great Lakes are highly uncertain (Gronewold and Rood, 2019),
237 in part because of uncertainty in future runoff draining into the lakes from a large contributing area
238 (Kayastha et al. 2022), much of which is ungauged (Fry et al., 2013). Improved rainfall-runoff models that
239 can regionalize across the entire Great Lakes basin are necessary to help address this challenge, and so an
240 auxiliary goal of this work is to contribute PIML rainfall-runoff models to the Great Lakes Runoff
241 Intercomparison Project Phase 4 (~~GRIP-GL~~) presented in Mai et al. (2022). This study currently provides
242 one of the most robust benchmarks comparing DL rainfall-runoff models to a range of process-based
243 models, and so we design our experiment to be consistent with the data and model development rules
244 outlined in ~~the GRIP-GL~~ [that intercomparison project](#).

246 **2. Data**

247 This study focuses on 212 watersheds draining into the Great Lakes and Ottawa River, which are all located
248 in the St. Lawrence River basin (Figure 1). ~~We note that this region is of similar spatial scale to other~~
249 ~~benchmarking datasets for DL rainfall-runoff models (e.g., CAMELS-GB; Coxon et al., 2020).~~ For direct
250 comparability to previous results from the [Great Lakes Runoff Intercomparison Project](#) ~~GRIP-GL~~, all data
251 for these watersheds are taken directly from the work in Mai et al. (2022) and include daily streamflow time
252 series, meteorological forcings, geophysical attributes for each watershed, and auxiliary hydrologic fluxes.
253 Daily streamflow were gathered from the U.S. Geological Survey (~~USGS~~) and Water Survey Canada (~~WSC~~)
254 between January 2000 and December 2017. All streamflow gauging stations have a drainage area greater
255 than or equal to 200 km² and less than 5% missing data in the study period. The watersheds are evenly
256 distributed across the five lake basins and the Ottawa River basin, and they represent a range of land

257 use/land cover types and degrees of hydrologic alteration from human activity. In the experiments described
 258 further below, 141 of the watersheds are designated as training sites, and the remaining 71 watersheds are
 259 used for testing (see Figure 1). In addition, the period between January 2000 to December 2010 is reserved
 260 for model training (termed the training period), and the period between January 2011 – December 2017 is
 261 used for model testing (termed the testing period).

262



263

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

267

268 Meteorological forcings are taken from the Regional Deterministic Reanalysis System v2 (~~RDRS-v2~~),
 269 which is an hourly, 10 km dataset available across North America (Gasset et al., 2021). Hourly precipitation,
 270 net incoming shortwave radiation (R_s), and temperature are aggregated into a basin-wide daily precipitation
 271 average, daily R_s average, and daily minimum and maximum temperature. We note that the precipitation
 272 data from [the Regional Deterministic Reanalysis System v2 ~~RDRS-v2~~](#) is produced from the Canadian
 273 Precipitation Analysis (~~CaPA~~), which combines available surface observations of precipitation with a short-

274 term reforecast provided by the 10 km Regional Deterministic Reforecast System. That is, the precipitation
 275 data is not model based, but rather is based on gauged data and spatially interpolated using information
 276 from modeled output.

277
 278 Geophysical attributes for each watershed were collected from a variety of sources. Basin-average statistics
 279 of elevation and slope were derived from the HydroSHEDS dataset (Lehner et al., 2008), which provides a
 280 digital elevation model (DEM) with 3 arcsec resolution. Soil properties (e.g., soil texture, classes) were
 281 gathered from the Global Soil Dataset for Earth System Models (GSDE; Shangguan et al., 2014), which is
 282 available at a 30 arcsec resolution. Land cover data at a 30 m resolution and based on Landsat imagery from
 283 2010-2011 were derived from the North American Land Change Monitoring System (NALCMS, 2017).
 284 These geophysical datasets were used to derive basin-averaged attributes for each watershed, listed in Table
 285 1.

286
 287 **Table 1.** Watershed attributes used in the deep learning models developed in this work (adapted from Mai
 288 et al., 2022).

Attribute	Description
p_mean	Mean daily precipitation
pet_mean	Mean daily potential evapotranspiration
aridity	Ratio of mean PET to mean precipitation
t_mean	Mean of daily maximum and daily minimum temperature
frac_snow	Fraction of precipitation falling on days with mean daily temperatures below 0°C
high_prec_freq	Fraction of high-precipitation days (= 5 times mean daily precipitation)
high_prec_dur	Average duration of high-precipitation events (number of consecutive days with = 5 times mean daily precipitation)
low_prec_freq	Fraction of dry days (< 1 mm d-1 daily precipitation)
low_prec_dur	Average duration of dry periods (number of consecutive days with daily precipitation < 1 mm d-1)
mean_elev	Catchment mean elevation
std_elev	Standard deviation of catchment elevation

mean_slope	Catchment mean slope
std_slope	Standard deviation of catchment slope
area_km2	Catchment area
Temperate-or-sub-polar-needleleaf-forest	Fraction of land covered by “Temperate-or-sub-polar-needleleaf-forest”
Temperate-or-sub-polar-grassland	Fraction of land covered by “Temperate-or-sub-polar-grassland”
Temperate-or-sub-polar-shrubland	Fraction of land covered by “Temperate-or-sub-polar-shrubland”
Temperate-or-sub-polar-grassland	Fraction of land covered by “Temperate-or-sub-polar-grassland”
Mixed-Forest	Fraction of land covered by “Mixed-Forest”
Wetland	Fraction of land covered by “Wetland”
Cropland	Fraction of land covered by “Cropland”
Barren-Lands	Fraction of land covered by “Barren-Lands”
Urban-and-Built-up	Fraction of land covered by “Urban-and-Built-up”
Water	Fraction of land covered by “Water”
BD	Soil bulk density (g cm ⁻³)
CLAY	Soil clay content (% of weight)
GRAV	Soil gravel content (% of volume)
OC	Soil organic carbon (% of weight)
SAND	Soil sand content (% of weight)
SILT	Soil silt content (% of weight)

289

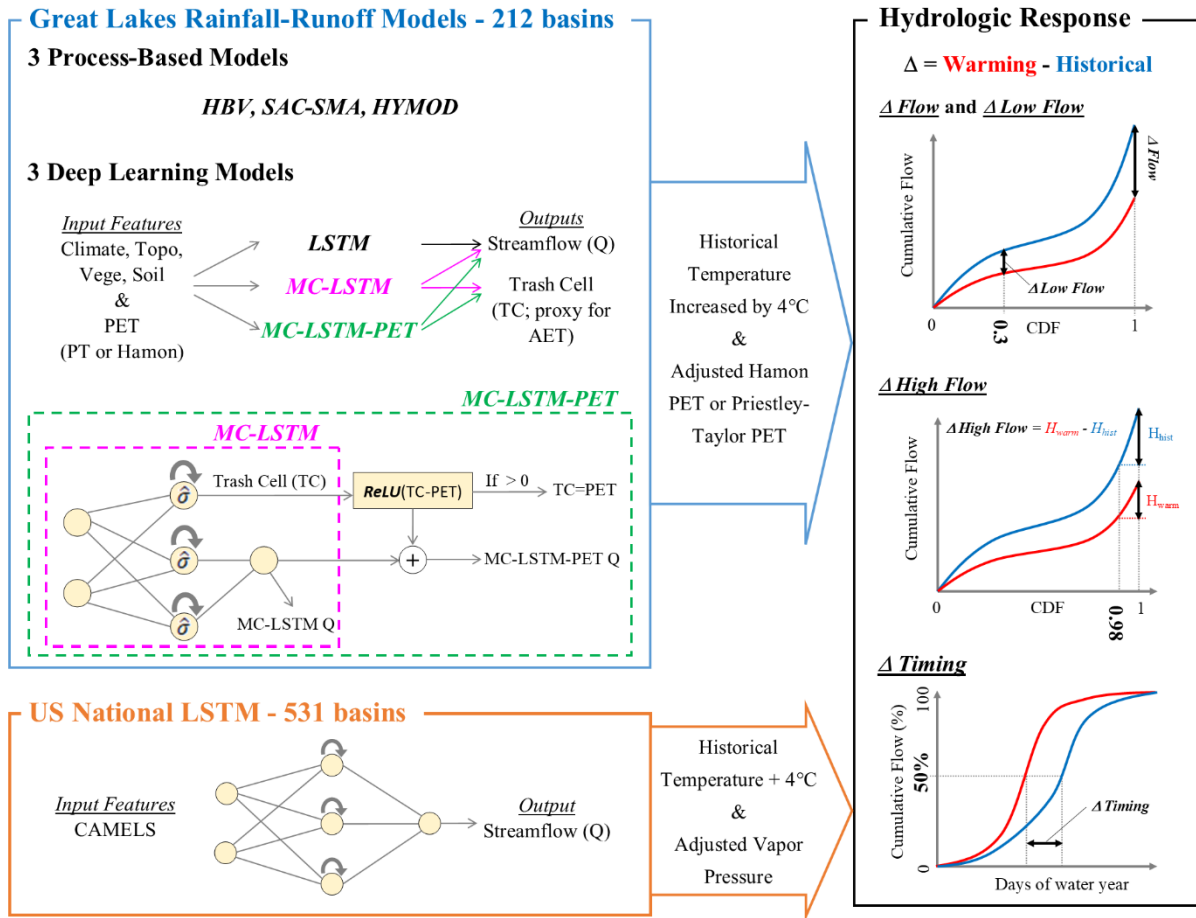
290 Finally, we also collect daily actual evapotranspiration (AET) for each watershed in millimeters per day,
291 which was originally taken from the Global Land Evaporation Amsterdam Model (GLEAM) v3.5b dataset
292 (Martens et al., 2017). GLEAM couples remotely sensed observations of microwave Vegetation Optical
293 Depth, a multi-layer soil moisture model driven by observed precipitation and assimilating satellite surface
294 soil moisture observations, and Priestly-Taylor based estimates of PET to derive an estimate of AET for
295 each day. The daily data were originally available over the entire study domain at a 0.25° resolution between
296 2003-2017 and were aggregated to basin-wide totals for each watershed. [While AET from GLEAM is still](#)
297 [uncertain, it provides a useful, independent, remote-sensing based benchmark against which to compare](#)
298 [rainfall-runoff model estimates of AET.](#)

299

300 3. Methods

301 We design an experiment to test the two primary hypotheses of this study, namely that a standard LSTM
302 will overestimate ~~hydrologic-water~~ losses under warming because of an overreliance on historical
303 correlations between temperature and PET, while this effect will be lower in PIML-based rainfall-runoff
304 models designed to better account for water loss in the system. To conduct this experiment, we develop
305 three different DL rainfall-runoff models to predict daily streamflow across the Great Lakes region, as well
306 as three conceptual, process-based models as benchmarks, each of which is trained twice with either an
307 energy budget-based or temperature-based estimate of PET. The DL models include a regional LSTM very
308 similar to the model in Mai et al., (2022), an MC-LSTM that conserves mass, and a new variant of the MC-
309 LSTM that also respects the relationship between PET and water loss (termed MC-LSTM-PET). After
310 comparing historical model performance, we ~~conduct a sensitivity analysis foree on~~ all models ~~with climate~~
311 ~~change scenarios in which composed of~~ historical ~~precipitation and historical but warmed~~ temperatures are
312 warmed by 4°C, as well as PET is updated based on those warmed temperatures, and all other
313 meteorological variable time series are left unchanged from historical values. This is a similar approach to
314 that taken in SW22, but in contrast to that study this work 1) focuses on the magnitude of streamflow
315 response to warming under two different PET formulations; 2) considers a different set of physics-informed
316 DL models in which the architecture (rather than the inputs or targets) of the model are changed to better
317 preserve physical plausibility under ~~unprecedented shifts in climate-change~~; and 3) evaluates an expanded
318 set of hydrologic metrics to better understand both the plausibility and the variability of ~~climate-change~~
319 responses across the different models. Finally, in a subset of the analysis, we also utilize a fourth DL model,
320 the LSTM used in SW22 that was previously fit to 531 basins across the CONUS (Kratzert et al. 2021),
321 which uses daily precipitation, maximum and minimum temperature, radiation, and vapor pressure as input
322 but not PET. This model is used to evaluate whether a DL model fit to many more watersheds that span a
323 more diverse gradient of climate conditions behaves differently under warming than an LSTM fit only to
324 locations in the Great Lakes basin. Figure 2 presents an overview of our experimental design.

325



327

328 **Figure 2.** Overview of experiment design. Three deep learning rainfall-runoff models (LSTM, MC-
 329 LSTM, MC-LSTM-PET) and three conceptual, process-based models (HBV, SAC-SMA, HYMOD) are
 330 trained and tested across 212 watersheds throughout the Great Lakes basin. Models are validated by
 331 comparing predictions to streamflow (Q) and actual evapotranspiration (AET). All models are then forced
 332 with historical meteorology, but with historical temperatures warmed by 4°C and potential
 333 evapotranspiration (PET) updated based on those warmed temperatures using either the Hamon or
 334 Priestley-Taylor method. Hydrologic model responses across all models are then compared in terms of
 335 long-term mean daily flows, low flows, high flows, and streamflow seasonal timing statistics. The
 336 experiment is also repeated with an LSTM fit to 531 basins across the contiguous United States, except
 337 that model does not use PET as an input and vapor pressure is also adjusted along with temperature.
 338

329 3.1. Models

340 3.1.1. Benchmark Conceptual Models

341 We develop three conceptual, process-based hydrologic models as benchmarks, including the Hydrologiska
 342 Byråns Vattenbalansavdelning (HBV) model (Bergström and Forsman, 1973), HYMOD (Boyle, 2001), and

343 the Sacramento Soil Moisture Accounting (SAC-SMA) model (Burnash, 1995) coupled with SNOW-17
344 (Anderson, 1976). These models are developed as lumped, conceptual models for each watershed, and were
345 selected for several reasons. First, in the Great Lakes Intercomparison Project (Mai et al., 2022), HYMOD
346 was one the best performing process models for both streamflow and AET estimation. SAC-SMA is widely
347 used in the United States, forming the core hydrologic model in NOAA’s Hydrologic Ensemble Forecasting
348 System (Demargne et al., 2014). We also found in WS22 that AET from SAC-SMA matched the seasonal
349 pattern of MODIS-derived AET well across California. HBV is also an extremely popular model (Seibert
350 and Bergström, 2022), is used for operational forecasting in multiple countries (Olsson and Lindstrom,
351 2008; Krøgli et al., 2018), and performs very well in hydrologic model intercomparison projects (Breuer et
352 al., 2009; Plesca et al., 2012; Beck et al., 2016, 2017).

353
354 We calibrate the process-based models with the genetic algorithm from Wang et al. (1991) to maximize
355 minimize the Nash-mean-Sutcliffe-squared Efficiency-error (NSEMSE), using a population size equal to
356 100 times the number of parameters, evolved over 100 generations, and with a spin-up period of 1 year.
357 Each benchmark model is calibrated separately to each of the 141 training sites using the temporal train/test
358 split described in Section 2, and training is repeated- 10 separate times with different random initializations
359 to account for uncertainty in the training process and to estimate parametric uncertainty. Benchmark models
360 are developed for the 71 testing sites in two ways: 1) separate models are trained for the testing sites during
361 the training period; and 2) each testing site is assigned a donor from among the 141 training sites, and the
362 calibrated parameters from that donor site are transferred to the testing site. The first of these approaches
363 enables a comparison between DL models fit only to the training sites to benchmark models developed for
364 the testing sites, i.e., a spatial out-of-sample versus in-sample comparison. The second of these approaches
365 enables a more direct spatial out-of-sample comparison between DL and benchmark models. We note that
366 donor sites were used to assign model parameters to testing sites in the benchmarking study of Mai et al.
367 (2022), and to retain direct comparability to the results of that work we use the same donor sites for each

368 testing site. Donor sites were selected based on spatial proximity, while also prioritizing donor sites that
 369 were nested within the watershed of the testing site.

370

371 3.1.2. LSTM

372 We develop a single, regional LSTM for predicting daily streamflow across the Great Lakes region. In the
 373 LSTM, nodes within hidden layers feature gates and cell states that address the vanishing gradient problem
 374 of classic recurrent neural networks and help capture long-term dependencies between input and output
 375 time series. The model defines a D -dimensional vector of recurrent cell states $\mathbf{c}[t]$ that is updated over a
 376 sequence of $t=1, \dots, T$ time steps based on a sequence of inputs $\mathbf{x} = \mathbf{x}[1], \dots, \mathbf{x}[T]$, where each input $\mathbf{x}[t]$ is
 377 a K -dimensional vector of features. Information stored in the cell states is then used to update a D -
 378 dimensional vector of hidden states $\mathbf{h}[t]$, which form the output of the hidden layer in the model. The
 379 structure of the LSTM is given as follows:

380

$$381 \quad \mathbf{i}[t] = \sigma(\mathbf{W}_i \mathbf{x}[t] + \mathbf{U}_i \mathbf{h}[t-1] + \mathbf{b}_i) \quad (\text{Eq. 1.1})$$

$$382 \quad \mathbf{f}[t] = \sigma(\mathbf{W}_f \mathbf{x}[t] + \mathbf{U}_f \mathbf{h}[t-1] + \mathbf{b}_f) \quad (\text{Eq. 1.2})$$

$$383 \quad \mathbf{g}[t] = \tanh(\mathbf{W}_g \mathbf{x}[t] + \mathbf{U}_g \mathbf{h}[t-1] + \mathbf{b}_g) \quad (\text{Eq. 1.3})$$

$$384 \quad \mathbf{o}[t] = \sigma(\mathbf{W}_o \mathbf{x}[t] + \mathbf{U}_o \mathbf{h}[t-1] + \mathbf{b}_o) \quad (\text{Eq. 1.4})$$

$$385 \quad \mathbf{c}[t] = \mathbf{f}[t] \odot \mathbf{c}[t-1] + \mathbf{i}[t] \odot \mathbf{g}[t] \quad (\text{Eq. 1.5})$$

$$386 \quad \mathbf{h}[t] = \mathbf{o}[t] \odot \tanh(\mathbf{c}[t]) \quad (\text{Eq. 1.6})$$

$$387 \quad \mathbf{y}[T] = \text{ReLU}(\mathbf{W}_y \mathbf{h}[T] + b_y) \quad (\text{Eq. 1.7})$$

388

389 Here, the input gate ($\mathbf{i}[t]$) controls how candidate information ($\mathbf{g}[t]$) from inputs and previous hidden states
 390 flows to the current cell state ($\mathbf{c}[t]$); the forget gate ($\mathbf{f}[t]$) enables removal of information within the cell
 391 state over time; and the output gate ($\mathbf{o}[t]$) controls information flow from the current cell state to the hidden
 392 layer output. All bolded terms are vectors, and \odot denotes element-wise multiplication. To produce

393 streamflow predictions, $\mathbf{h}[T]$ at the last time step in the sequence is passed through a fully connected layer
 394 to a single-node output layer (i.e., a many-to-one formulation). We ensure nonnegative streamflow
 395 predictions using the rectified linear unit (ReLU) activation function for the output neuron, expressed as
 396 $\text{ReLU}(x) = \max(0, x)$. Importantly, there are no constraints requiring the mass of water entering as
 397 precipitation to be conserved within this architecture.

398
 399 The LSTM takes $K=39$ input features: 9 dynamic and 30 static. The dynamic input features are basin-
 400 averaged climate, including daily precipitation, maximum temperature, minimum temperature, net
 401 incoming shortwave radiation, specific humidity, surface air pressure, zonal and meridional components of
 402 wind, and PET. The static features represent catchment attributes (see Table 1) and are repeated for all time
 403 steps in the input sequences \mathbf{x} . All input features are standardized before training (by subtracting the mean
 404 and dividing by the standard deviation for data across all training sites in the training period). Note that we
 405 do not standardize the observed streamflow, besides dividing [my-by](#) drainage area to represent streamflow
 406 in units of millimeters.

407
 408 We train the LSTM by minimizing the mean-squared error averaged over the 141 training watersheds
 409 during the training period:

$$410 \quad \text{MSE} = \frac{1}{N} \sum_{n=1}^N \frac{1}{T_n} \sum_{t=1}^{T_n} (\hat{Q}_{n,t} - Q_{n,t})^2 \quad (2)$$

411 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
 412 $Q_{n,t}$ are, respectively, the streamflow prediction and observation for basin n and day t . To estimate $\hat{Q}_{n,t}$,
 413 we feed into the network an input sequence for the past $T=365$ days. The model was developed with 1
 414 hidden layer composed of $D=256$ nodes, a mini-batch size of 256, a learning rate of 0.0005, and a drop-out
 415 rate of 0.4, and it was trained across 30 epochs. All hyperparameters (number of hidden layer nodes, mini-
 416 batch size, learning rate, dropout rate, and number of epochs) were selected in a 5-fold cross-validation on
 417 the training sites ([see Table S2 for details on grid search](#)). Network weights are tuned using the ADAM

418 optimizer (Kingma & Ba, 2015). The model is trained 10 separate times with different random
 419 initializations to account for uncertainty in the training process.

420

421 For the evaluation of streamflow [projections-responses to under- climate-change warming](#), we also use an
 422 LSTM taken from Kratzert et al. (2021) and employed in SW22, which was fit to 531 basins across the
 423 contiguous United States (hereafter called the National LSTM). This model was trained using a different
 424 set of data compared to our Great Lakes LSTM but also used a mix of dynamic and static features, all of
 425 which were drawn from the Catchment Attributes and Meteorology for Large-Sample Studies (CAMELS)
 426 dataset (Newman et al., 2015). This model uses daily precipitation, maximum and minimum temperature,
 427 shortwave downward radiation, and vapor pressure as input but not PET. However, we note that
 428 temperature, radiation, and vapor pressure are the three major inputs (besides wind speeds) needed to
 429 calculate energy budget-based PET. [There are 29 CAMELS watersheds located within the Great Lakes](#)
 430 [basin, and 17 of those 29 watersheds were also used in the training and testing sets for the Great Lakes](#)
 431 [LSTM \(see Figure 1\).](#)

432

433 3.1.3. MC-LSTM

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

438

$$439 \quad \hat{\mathbf{c}}[t - 1] = \frac{\mathbf{c}[t-1]}{\|\mathbf{c}[t-1]\|_1} \quad (\text{Eq. 3.1})$$

$$440 \quad \mathbf{i}[t] = \hat{\sigma}(\mathbf{W}_i \mathbf{x}[t] + \mathbf{U}_i \hat{\mathbf{c}}[t - 1] + \mathbf{V}_i \mathbf{a}[t] + \mathbf{b}_i) \quad (\text{Eq. 3.2})$$

$$441 \quad \mathbf{o}[t] = \sigma(\mathbf{W}_o \mathbf{x}[t] + \mathbf{U}_o \hat{\mathbf{c}}[t - 1] + \mathbf{V}_o \mathbf{a}[t] + \mathbf{b}_o) \quad (\text{Eq. 3.3})$$

$$442 \quad \mathbf{R}[t] = \hat{\sigma}(\mathbf{W}_R \mathbf{x}[t] + \mathbf{U}_R \hat{\mathbf{c}}[t - 1] + \mathbf{V}_R \mathbf{a}[t] + \mathbf{b}_R) \quad (\text{Eq. 3.4})$$

443 $\mathbf{m}[t] = \mathbf{R}[t]\mathbf{c}[t - 1] + \mathbf{i}[t]\mathbf{x}[t]$ (Eq. 3.5)

444 $\mathbf{c}[t] = (1 - \mathbf{o}[t]) \odot \mathbf{m}[t]$ (Eq. 3.6)

445 $\mathbf{h}[t] = \mathbf{o}[t] \odot \mathbf{m}[t]$ (Eq. 3.7)

446

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

453

454 $\hat{\sigma}(z_j) = \frac{\sigma(z_j)}{\sum_j \sigma(z_j)}$ (Eq. 4)

455

456 Here, $\sigma(\cdot)$ is the sigmoid activation function, while $\hat{\sigma}(\cdot)$ is a normalized sigmoid activation that produces a
 457 vector of fractions that sum to unity.

458

459 The mass in $\mathbf{m}[t]$, which is stored across D elements in the vector, is then distributed to the output of the
 460 hidden layer, $\mathbf{h}[t]$, or the next cell state, $\mathbf{c}[t]$. To account for water losses from evapotranspiration or other
 461 sinks, one element of the D -dimensional vector $\mathbf{h}[t]$ is considered a ‘trash cell’, and the output of this cell
 462 is ignored when calculating the final streamflow prediction, which at time T is given by the sum of outgoing
 463 water mass:

464

465 $y[T] = \sum_{d=1}^{D-1} h_d[T]$ (Eq. 5)

466

467 Here, the D^{th} cell of \mathbf{h} (h_D) is set as the trash cell, and water allocated to this cell at each time step $t=1,\dots,T$
468 is lost from the system. We note that the MC-LSTM was trained in the same way as the LSTM (i.e., same
469 inputs, loss function, training and test sets, hyperparameter selection process, number of ensemble members
470 with random initialization).

471

472 **3.1.4. MC-LSTM-PET**

473 We also propose a novel variant of the MC-LSTM that requires water lost from the system to not exceed
474 PET (hereafter referred to as the MC-LSTM-PET). In the original MC-LSTM, any amount of water can be
475 delegated to the trash cell h_D . Therefore, while water is conserved in the MC-LSTM, the model has the
476 freedom to transfer any amount of water from $\mathbf{m}[t]$ to the trash cell (and out of the hydrologic system) as
477 it seeks to improve the loss function during training. This has the benefit of handling biased data, e.g., cases
478 where the precipitation input to the system is systematically too high compared to the measured outflow.
479 However, this structure also has the drawback of potentially removing more water from the system than is
480 physically plausible. To address this issue, we propose a small change to the architecture of the MC-LSTM,
481 where any water relegated to the trash cell that exceeds PET at time t is directed back to the stream:

482

$$483 \quad y[t] = \sum_{d=1}^{D-1} h_d[t] + \text{ReLU}(h_D[t] - \text{PET}[t]) \quad (\text{Eq. 6})$$

484

485 Here, the ReLU activation ensures that any water in the trash cell (h_D) which exceeds PET at time t is
486 added to the streamflow prediction $y[t]$, but the streamflow prediction is the same as the original MC-
487 LSTM (Eq. 5) if water in the trash cell is less than PET. [This approach assumes that the maximum allowable
488 water lost from the system cannot exceed PET, and therefore ignores other potential terminal sinks \(e.g.,
489 inter-basin lateral groundwater flows; human diversions and inter-basin transfers\). This assumption is more
490 strongly supported in moderately-sized \(\$> 200 \text{ km}^2\$ \), low-gradient, non-arid watersheds where inter-basin
491 groundwater flows are less impactful \(Fan 2019; Gordon et al., 2022\), such as the Great Lakes basins](#)

492 ~~examined in this work. However, we discuss the potential to relax the assumptions of the MC-LSTM-PET~~
493 ~~model in Section 5. This approach assumes that the maximum allowable water lost from the system cannot~~
494 ~~exceed PET, and therefore ignores other potential terminal sinks (e.g., deep groundwater percolation that~~
495 ~~remains disconnected from the stream; lateral groundwater flows out of the watershed; human diversions).~~
496 ~~However, given that evapotranspiration accounts for the vast majority of water lost in most hydrologic~~
497 ~~systems, this assumption is likely reasonable in most cases.~~ The MC-LSTM-PET was trained in the same
498 way as the LSTM (i.e., same inputs, loss function, training and test sets, hyperparameter selection process,
499 number of ensemble members with random initialization).

500

501 **3.2. Model Performance Evaluation**

502 As noted previously, 141 of the watersheds are designated as training sites, and the remaining 71 watersheds
503 are used for testing. In addition, the training and testing periods were restricted to January 2000 -December
504 2010 and January 2011 – December 2017, respectively. This provides three separate ways to evaluate model
505 performance:

- 506 • Temporal validation - Performance across models is evaluated at training sites during the testing
507 period.
- 508 • Spatial validation - Performance across models is evaluated at testing sites during the training
509 period.
- 510 • Spatiotemporal validation - Performance across models is evaluated at testing sites during the
511 testing period.

512

513 All three evaluation strategies are utilized. For benchmark process-based models that are calibrated locally
514 on a site-by-site basis, we consider model versions that are transferred to testing sites from training sites,
515 as well as models that are trained to the testing sites directly (see Section 3.1.1). The former can be used

516 for all three evaluation strategies above, while the latter can only be used for temporal validation at the
517 testing sites.

518
519 [Following other intercomparison studies \(Frame et al., 2022; Gauch et al., 2021a; Klotz et al., 2022; Kratzert](#)
520 [et al., 2021\)](#), ~~Several-several~~ metrics are considered for model evaluation, including percent bias (PBIAS),
521 the Nash-Sutcliffe Efficiency (NSE; Nash and Sutcliffe, 1970), Kling-Gupta Efficient (KGE; Gupta et al.
522 2009), top 2% peak flow bias (FHV; Yilmaz et al. 2008), and bottom 30% low flow bias (FLV; Yilmaz et
523 al. 2008). Each metric is calculated separately for training and testing periods for each site. For ~~all the DL~~
524 models, all results are estimated from the ensemble mean from 10 separate training trials.

525
526 For the process models, the MC-LSTM, and the MC-LSTM-PET, we also compare simulations of AET to
527 ~~observations of~~ AET from the GLEAM database. We note that AET data were not used to train any of the
528 models. For the process models, AET is a direct output of the model and so can immediately be extracted
529 for comparison, but AET is not directly simulated by the MC-LSTM or MC-LSTM-PET. Instead, we
530 assume water delegated to the trash cell permanently leaves the system because of evapotranspiration.
531 Several metrics are used to compare model-based AET to GLEAM AET, including KGE, correlation, and
532 PBIAS, and the comparison is conducted for training sites during the training period and under temporal,
533 spatial, and spatiotemporal validation (as described above). Similar to streamflow, all AET results ~~for the~~
534 ~~MC-LSTM and MC-LSTM-PET~~ are based on the ensemble mean ~~of water delegated to the trash cell~~ from
535 the 10 separate training trials.

536
537 **3.3. Evaluating Hydrologic Response under Warming**

538 All Great Lakes models in this study are trained twice with different PET estimates as input, including the
539 Hamon method (a temperature-based approach; Hamon, 1963) and the Priestley-Taylor method (an energy
540 budget-based approach; Priestley and Taylor, 1972). [We select the Hamon method because of its stronger](#)
541 [dependence on temperature compared to other temperature-based approaches that also depend on radiation](#)

542 [\(e.g., Hargreaves and Samani, 1985; Oudin et al., 2005\)](#). We select the Priestley-Taylor method based on
 543 [its widespread use in the literature \(Wu et al., 2021; Su and Singh, 2023\)](#) and its approximation of the more
 544 [physically-based Penman-Monteith approach \(Allen et al. 1998\)](#). Together, these two approaches lie
 545 [towards the lower and upper bounds of temperature sensitivity across multiple PET approaches \(see Shaw](#)
 546 [and Riha, 2011\)](#).

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

549

$$550 \quad PET_H = \alpha_H \times 29.8 \times Hr_{day} \frac{e_{sat}}{T_a + 273.2} \quad (\text{Eq. 7})$$

$$551 \quad e_{sat} = 0.611 \times \exp\left(\frac{17.27 \times T_a}{237.3 + T_a}\right) \quad (\text{Eq. 8})$$

552 where Hr_{day} is the number of daylight hours, T_a is the average daily temperature ($^{\circ}\text{C}$) calculated from
 553 daily minimum and maximum temperature, e_{sat} is the saturation vapor pressure (kPa), and α_H is a
 554 calibration coefficient set to 1.2 for all models in this study (similar to Lu et al., 2005).

555

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

557

$$558 \quad PET_{PT} = \alpha_{PT} \left(\frac{\Delta(T_a) \times (R_n - G)}{\lambda(\Delta(T_a) + \gamma)} \right) \times 1000 \quad (\text{Eq. 9})$$

559

560 Here, $\Delta(T_a)$ is the slope of the saturation vapor pressure temperature curve ($\text{kPa}/^{\circ}\text{C}$) and is a function of
 561 T_a , γ is the psychrometric constant ($\text{kPa}/^{\circ}\text{C}$), λ is the volumetric latent heat of vaporization (MJ/m^3), R_n is
 562 the net radiation ($\text{MJ}/\text{m}^2\text{-day}$) equal to the difference between net incoming shortwave (R_{nS}) and net
 563 outgoing longwave (R_{nL}) radiation, G is the heat flux to the ground ($\text{MJ}/\text{m}^2\text{-day}$), and α_{PT} is a dimensionless
 564 coefficient set to 1.1 for all models in this study (similar to Szilagyi et al., 2017). Details on how to calculate
 565 γ , $\Delta(T_a)$, and R_{nL} are available in Allen et al. (1998), and we assume $G=0$. Net shortwave radiation is given

566 by $R_{ns} = (1 - \zeta)R_s$, with $\zeta = .23$ the assumed albedo and R_s the incoming shortwave radiation. We note
567 that net outgoing longwave radiation R_{nl} is a function of maximum and minimum temperature, actual vapor
568 pressure, and R_s (see Eq. 39 in Allen et al. 1998). All exogenous meteorological inputs for the two methods
569 are derived from the [Regional Deterministic Reanalysis System v2 RDRS v2](#) (see Section 2). We note that
570 using $\alpha_H = 1.2$ and $\alpha_{PT} = 1.1$ leads to very similar [long-term average](#) PET estimates between the Hamon
571 and Priestley-Taylor methods under baseline climate conditions, helping to ensure their comparability. [We](#)
572 [also note that both PET series are highly correlated with daily average temperatures \(average Pearson](#)
573 [correlations across sites of 0.94 and 0.83 for Hamon and Priestley-Taylor PET, respectively\).](#)

574

575 We then ~~develop a simple climate change scenario~~ [conduct a sensitivity analysis of model response](#) in which
576 the historical minimum and maximum temperature time series are increased uniformly by 4 °C, and the two
577 PET estimates are updated using these warmed temperatures. We focus the ~~climate change~~-assessment on
578 training period data at the training sites, so that any differences in ~~climate change projections~~ [responses](#) that
579 emerge between the DL and process models are due to model structural differences and not the effects of
580 spatiotemporal regionalization. In the Priestley-Taylor method, we maintain historical values for R_s to isolate
581 how changes in temperature and its effect on $\Delta(T_a)$ and R_{nl} influence changes in PET. The use of historical
582 R_s is supported by the results from CMIP5 projections presented in Lai et al. (2022), but this assumption is
583 discussed further in Section 5.

584

585 We also ~~develop~~ [conduct](#) a similar ~~climate change scenario~~ [sensitivity analysis for on](#) the National LSTM,
586 which uses five dynamic input features from the CAMELS dataset (daily precipitation, maximum
587 temperature, minimum temperature, R_s , and water vapor pressure). Here, temperatures are warmed by 4°C,
588 while precipitation and R_s are held at historical values. There is a strong correlation between vapor pressure
589 and minimum temperature in the CAMELS dataset, since minimum temperature is used to estimate the
590 water vapor pressure (Newman et al., 2015). Thus, to run the National LSTM under warming, we also

591 adjust the vapor pressure input based on the change imposed to minimum temperature. This procedure is
 592 detailed in SW22.

593
 594 For both the Great Lakes DL models and the National LSTM, the dynamic inputs are adjusted based on the
 595 warming scenarios above. We also consider changes to ~~some of~~ the static input features that depend on
 596 temperature and PET in their calculation (e.g., pet_mean, aridity, t_mean, frac_snow; see [Table 1 for feature](#)
 597 [descriptions and Table 1-Supporting Information S1 and Table S1 for details on adjustments to these](#)
 598 [features](#)), and then run all models using two settings: 1) with ~~climate~~ changes only to the dynamic features,
 599 and 2) with ~~climate~~ changes to both dynamic [features](#) and to static features that depend on those dynamic
 600 [features](#). In total, there are six scenarios run in this work, which are shown in Table 2.

601
 602 **Table 2.** Overview of the setup for the different scenarios run in this analysis. All models are driven with
 603 temperatures warmed by 4°C. The Great Lakes models include the HBV, SAC-SMA, HYMOD, LSTM,
 604 MC-LSTM, and MC-LSTM models that are trained and tested to the 212 sites across the Great Lakes basin.
 605

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

606
 607
 608
 609
 610 Ultimately, for each model we compare hydrologic ~~projections-responses~~ under the warmed scenario to
 611 their values under the baseline scenario with no warming. For the National LSTM, we only consider basins
 612 in the CAMELS dataset within the Great Lakes Basin. For the process models, we also evaluate the
 613 uncertainty in hydrologic response based on the range predicted across the 10 different training trials, as a

614 [simple means to evaluate how parametric uncertainty influences the predictions.](#) We examine four different
615 metrics for this comparison, including:

- 616 • AVG.Q: the [long-term average-mean of daily streamflow runoff](#) across the entire series.
- 617 • FHV: the average of the top 2% peak flows.
- 618 • FLV: the average of the bottom 30% low flows.
- 619 • COM: the median center of mass across all [water](#) years, where the center of mass is defined as the
620 day of the [water](#) year by which half of the total annual flow has passed.

621

622 If our hypothesis is correct that the LSTM cannot distinguish water loss differences with different PET
623 [projections-series](#) but similar warming while process-based and PIML models can, we would expect that
624 under the LSTM using both PET [projectionsseries](#), [average-long-term mean](#) flow will decline [significantly](#)
625 [substantially](#) and with similar magnitude to the process models using the temperature-based PET method
626 but not the energy budget-based PET method. We would also expect the National LSTM to exhibit similar
627 behavior, even though it was able to learn from a larger set of watersheds across a more diverse range of
628 climate conditions. Finally, if our hypothesis is correct, we would expect the PIML models (MC-LSTM,
629 MC-LSTM-PET) to follow the process model [projections-responses](#) more closely across the two different
630 PET [projectionsseries](#), at least in terms of the difference in magnitude of [average-long-term mean](#)
631 streamflow declines. ~~For~~ [To facilitate a broader comparison](#) [inter-model comparison of DL and process-](#)
632 [based models under warming \(which is largely absent from the literature\)](#), we also explore the differences
633 in low flow (FLV), high flow (FHV), and [seasonal](#) timing (COM) metrics across all model versions, where
634 we have less reason to anticipate how DL and process models will differ in their [projections-responses](#) and
635 across PET formulations. [However, for responses like seasonal streamflow timing \(COM\), we do anticipate](#)
636 [that realistic responses should show a shift towards more streamflow earlier in the year, as warmer](#)
637 [temperatures lead to more precipitation falling as rain rather than snow and drive snowmelt earlier in the](#)
638 [spring.](#)

639

640 4. Results

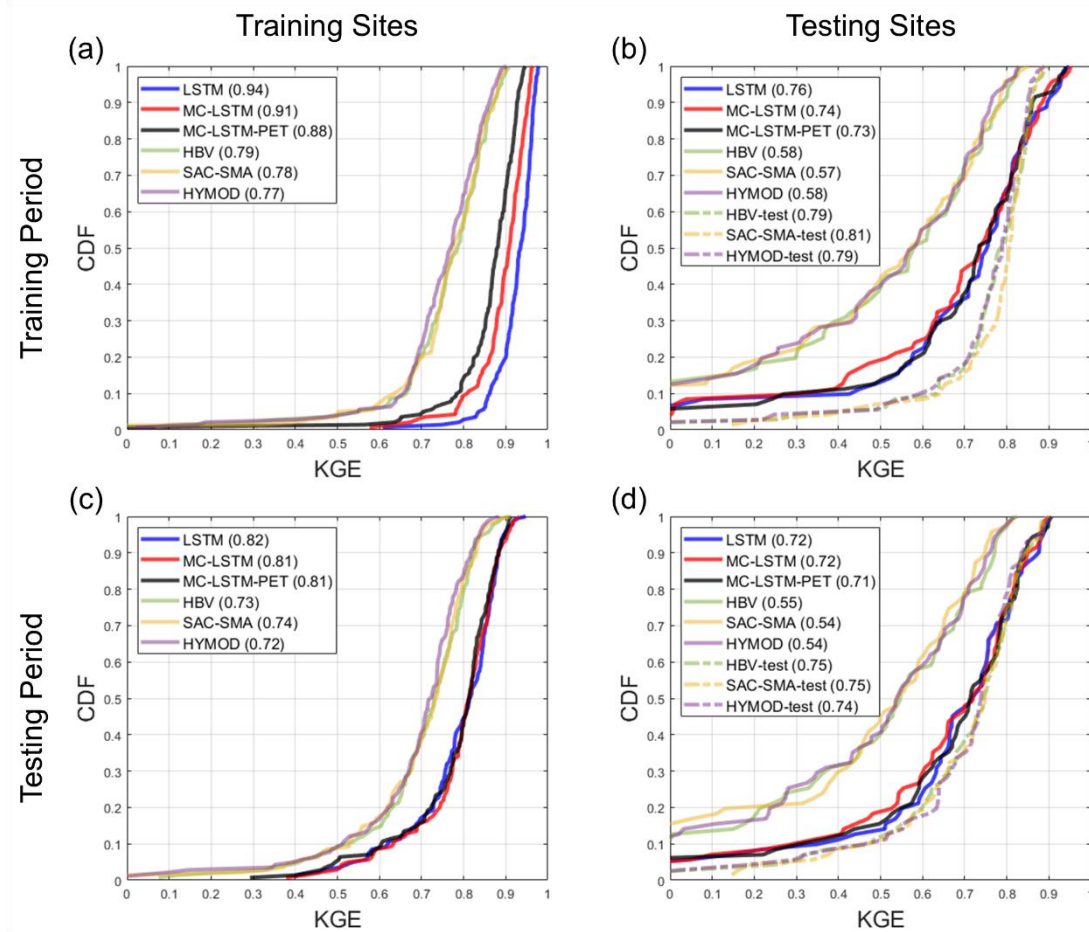
641 4.1. Model Performance Evaluation

642 Figure 3 shows the distribution of KGE values across sites for streamflow from the LSTM, MC-LSTM,
643 MC-LSTM-PET, and the three process-based models for both the training and testing sites during both the
644 training and testing periods. All results here and elsewhere in Section 4.1 are shown for the models fit with
645 Priestley-Taylor PET, but there is little difference in performance for the models fit with Hamon PET (see
646 Figure S1). For the process-based models, we show results for models fit to the training sites and then used
647 as donors at the testing sites, as well as models fit to the testing sites directly. We denote the latter with the
648 suffix “-test” and note that performance metrics at the training sites are not available for process models fit
649 to the testing sites.

650

651 Several insights emerge from Figure 3. First, for the training sites during the training period, all models
652 perform very well (Figure 3a). Across the three process models, the median KGE is ~~0.82~~0.79, ~~0.83~~0.78,
653 and ~~0.81~~0.77 for HBV, SAC-SMA, and HYMOD, respectfully. However, unsurprisingly, the DL models
654 perform better for the training data, with median KGE values all equal or above 0.88. The LSTM performs
655 best in this case. Under temporal validation (training sites during the testing period), performance degrades
656 somewhat across all models, and the differences in KGE between all process-based models and between
657 all DL models shrink considerably (Figure 3c). Larger performance declines are seen at the testing sites
658 during the training period (Figure 3b) and testing period (Figure 3d). Here, the median KGE for all process
659 models falls to between ~~0.56~~0.54-~~0.58~~0.57 when streamflow at the testing sites is estimated with donor models
660 from nearby gauged watersheds. In contrast, process models fit to the testing sites (denoted “-test”) exhibit
661 performance similar to that seen in Figure 3a,c. All three DL models perform quite well for the testing sites,
662 with median KGE values above 0.71 in both time periods. This is only modestly below the median KGE
663 for the process models fit to the testing sites, which is quite impressive given that this represents the spatial

664 out-of-sample performance of the DL models. We even see that for approximately 40% of testing sites
 665 during the training period, the DL models outperform the process models fit to those locations in that period.
 666



667
 668 **Figure 3.** The distribution of Kling-Gupta efficiency (KGE) for streamflow estimates across sites from
 669 each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results
 670 for the testing period are shown in panels (c) and (d), respectively. For the process models fit to the
 671 testing sites (denoted “-test”), no performance results are available at the training sites. All models are
 672 trained using Priestley-Taylor PET.
 673

674 Table 32 shows the median KGE, NSE, PBIAS, FHV, and FHL across testing sites for all models, excluding
 675 the process models fit to the testing sites. Similar to Figure 3, all three DL models outperform the donor-
 676 based process models at the testing sites for all metrics, ~~with the exception of PBIAS during the training~~
 677 ~~period~~. The performance across the three different DL models is similar, although there are some notable
 678 differences. In particular, the LSTM outperforms the MC-LSTM and MC-LSTM-PET for ~~KGE, NSE and,~~

679 ~~and~~ FLV ([as well as KGE in the training period](#)), the MC-LSTM-PET outperforms the LSTM and MC-
680 LSTM for PBIAS, and either the MC-LSTM or MC-LSTM-PET are the best performers for FHV. [The fact](#)
681 [that the MC-LSTM-PET performs best for PBIAS of all models suggests that the PET constraint imposed](#)
682 [in that model improves the overall accounting of water entering and existing the watershed on a long-term](#)
683 [basis](#). We [also](#) note that percent biases for FLV are high because the absolute magnitude of low flows is
684 small, so small absolute biases still lead to large percent biases.

685

686 **Table 32.** The median KGE, NSE, PBIAS, FHV, and FLV for streamflow across testing sites for the
687 training and testing periods for all models (excluding the process models fit to the testing sites). The metric
688 from the best performing model in each period is bolded. All models are trained using Priestley-Taylor PET.

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

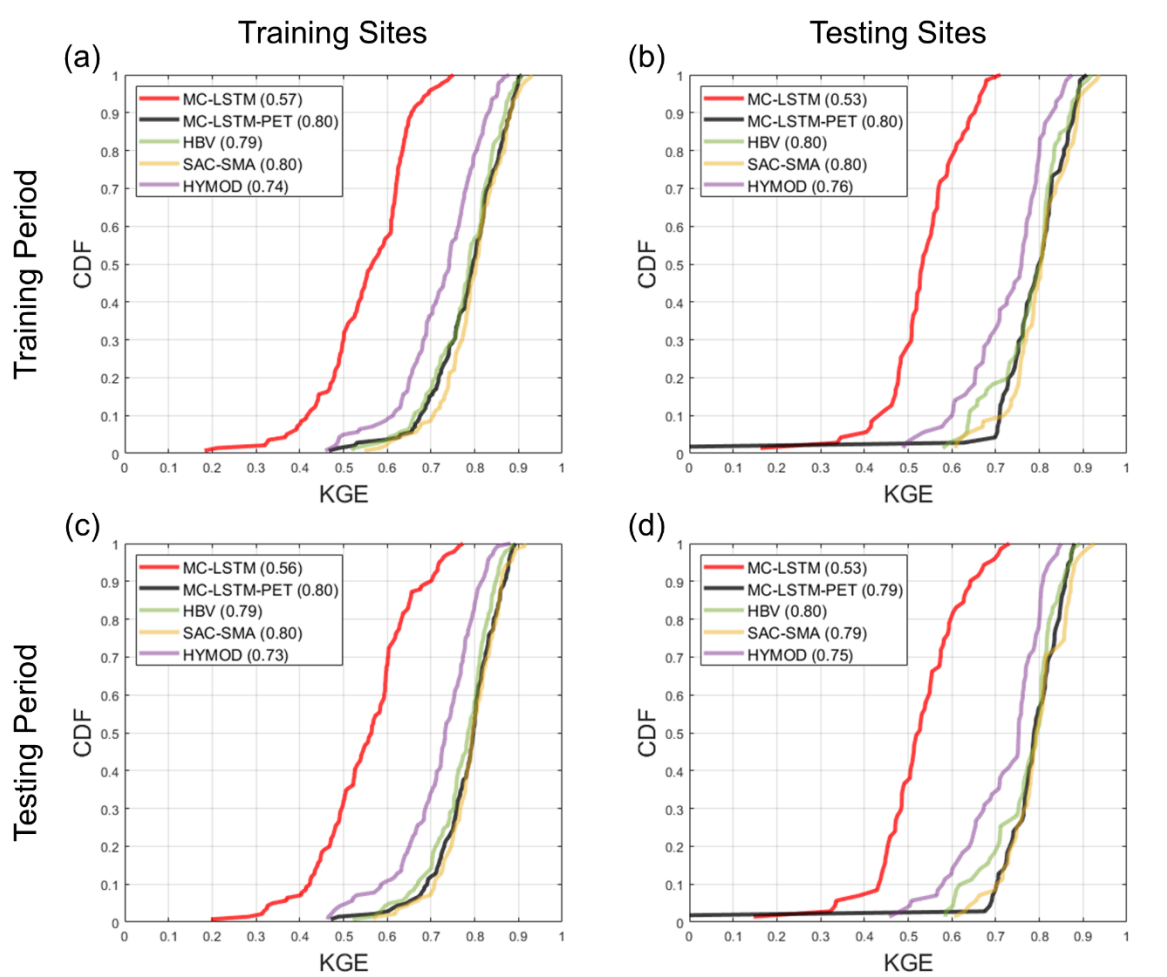
689

690 Figure 4 shows similar results as Figure 3, but for the KGE based on estimates of AET. Also, only donor
691 process models are shown for the testing sites. Results for correlation and PBIAS are available in the
692 Supplemental Information (Figures S2-S3). Here, the LSTM is not included because estimates of AET are
693 unavailable, while AET from the MC-LSTM and MC-LSTM-PET is based on water relegated to the trash
694 cell. Note that none of the models were trained for AET, and so results at training sites during the training
695 period also provide a form of model validation. Figure 4 shows that SAC-SMA and HBV predict AET with
696 relatively high degrees of accuracy for both training and testing sites in both periods (median KGE between
697 0.799-0.80). Performance is slightly worse for HYMOD. Notably, the MC-LSTM-PET exhibits very
698 similar, strong performance for all sites and periods as compared to SAC-SMA and HBV, except for one

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

701

702



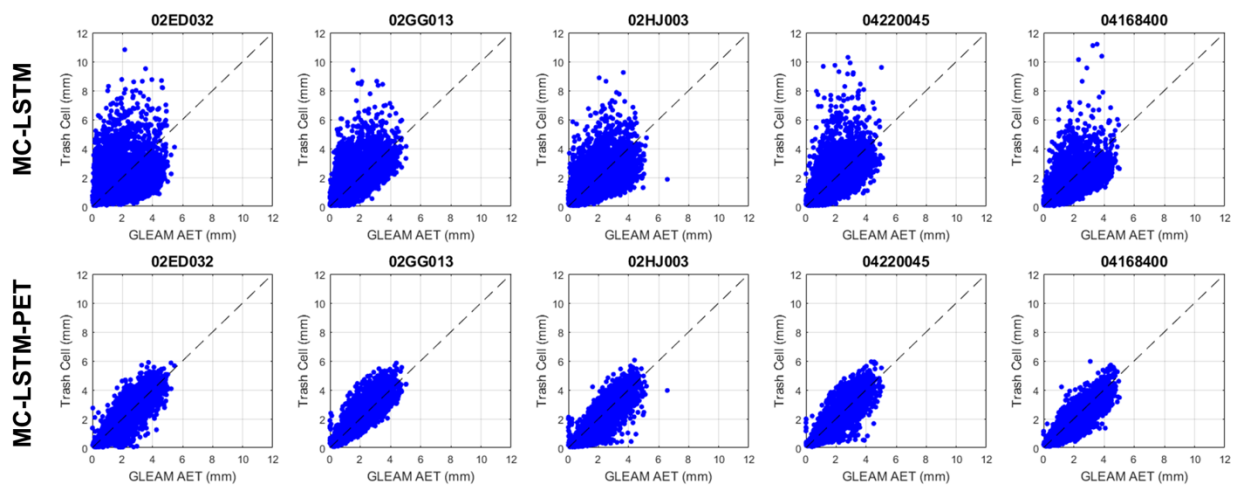
703

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

708

709 Further investigation reveals that the differences in KGE between the MC-LSTM and MC-LSTM-PET
710 models for AET are largely driven by differences in correlation (see Figure S2). We examine this difference
711 in more detail in Figure 5, which presents scatterplots of [observed-GLEAM](#) AET versus water allocations
712 to the trash cell for the two models from five randomly sampled testing sites across both training and testing

713 periods ([see Table S1 for details on each site](#)[Figure 1](#); also [Table S3](#)). Trash cell water from the MC-LSTM
 714 is not only more scattered around [observed-GLEAM AET](#) compared to the MC-LSTM-PET, but it also
 715 exhibits many outlier values that are two to five times larger than [observed-GLEAM AET](#). The MC-LSTM-
 716 PET follows the variability of [GLEAM](#) AET much more closely, with virtually no outliers that exceed
 717 [GLEAM](#) AET by large margins. This suggests that the PET constraint on the trash cell in the MC-LSTM-
 718 PET helps water allocated to that cell more faithfully represent [an-ET-sink](#)[evaporative water loss](#) in the DL
 719 model.



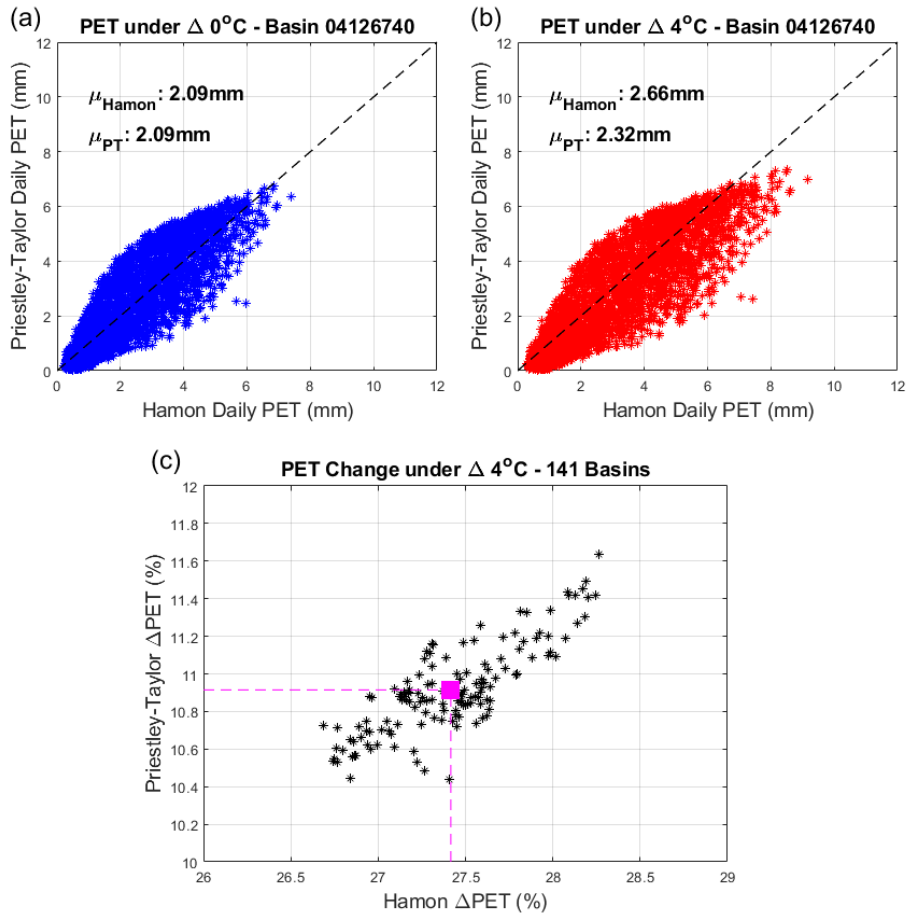
720
 721 **Figure 5.** Scatterplots of daily AET versus trash cell water for the (top) MC-LSTM and (bottom) MC-
 722 LSTM-PET at five randomly selected testing sites across both training and testing periods. All models are
 723 trained using Priestley-Taylor PET.
 724

725 **4.2. Evaluating Hydrologic Response under Warming**

726 Next, we evaluate streamflow [projections-responses](#) under a 4 °C warming scenario. We focus on training
 727 sites during the training period, so that any differences that emerge between DL and process models are
 728 only related to model structure and not spatiotemporal regionalization. [However, our results are largely](#)
 729 [unchanged if based on responses for testing sites in the testing period](#) ([see Figure S4](#)). First, we show the
 730 differences in historic and [warming-projected-adjusted](#) PET when using the Hamon and Priestley-Taylor
 731 methods ([Figure 6](#)). For the training period without any temperature change, PET estimated from the two
 732 methods is very similar ([Figure 6a](#); shown at one sample location for demonstration, [see Table S1](#)[Figure 1](#)

733 [and Table S3; Figure 6a](#)). However, under the scenario with 4 °C of warming, Hamon-based PET is
 734 [significantly-substantially](#) larger than Priestley-Taylor based PET (Figure 6b). On average, this difference
 735 reaches ~16% across all training sites and exhibits very little variability across locations (Figure 6c). The
 736 primary reason for the difference in [projected-the estimated](#) change in PET is that the Hamon method
 737 attributes PET entirely to temperature, while only a portion of PET is based on temperature in the Priestley-
 738 Taylor method, with the rest based on R_n . It is worthwhile to note that R_n does [change-increase](#) with
 739 temperature through its effects on net outgoing longwave radiation, but these changes [are-small-are generally](#)
 740 [less than 5% across all sites \(Allen et al. 1998\)](#).

741



742

743 **Figure 6.** (a) Daily PET estimated using the Hamon and Priestley-Taylor method for one sample
 744 watershed, under historic climate conditions in the training period. (b) Same as (a), but under the [climate](#)

745 ~~change~~ scenario with 4 °C of warming. (c) Percent change in average PET with 4 °C of warming across
746 all training sites using the Hamon and Priestley-Taylor methods.
747

748 Figure 7 shows how these differences in PET under warming propagate into changes in different attributes
749 of streamflow across training sites in the training period. The left and right columns of Figure 7 show
750 ~~projections-streamflow responses~~ using Hamon and Priestley-Taylor PET, respectively, while the rows of
751 Figure 7 show the distribution of changes (~~as a percentage~~) in different streamflow attributes (AVG.Q, FLV,
752 FHV, COM) across models. Figure 7 shows results for DL models where only the dynamic inputs are
753 changed under warming, ~~while Figure S4 show the same results when both the dynamic and the static~~
754 ~~climate properties are updated with warming.~~
755

756 Starting with changes in AVG.Q, Figure 7a,b shows that under the Hamon method for PET, the DL models
757 exhibit similar changes in ~~average-long-term mean~~ streamflow to the process-based models, with the
758 median Δ AVG.Q across sites ranging between -17% and ~~-2325%~~ across all models. However, when using
759 Priestley-Taylor PET, larger differences in the distribution of Δ AVG.Q emerge. Across all three process
760 models, the median Δ AVG.Q is between ~~-56%~~ to ~~-409%~~, and very few locations exhibit Δ AVG.Q less than
761 -20%. Conversely, the LSTM shows a median water loss of -20% under Priestley-Taylor PET and a very
762 similar distribution of water losses regardless of whether Hamon or Priestley-Taylor PET was used. The
763 MC-LSTM is also relatively insensitive to PET, and as compared to the process models, the MC-LSTM
764 tends to predict smaller absolute changes to AVG.Q for Hamon PET and larger changes under Priestley-
765 Taylor PET. Only the MC-LSTM-PET model achieves water loss that is ~~significantly-considerably~~ smaller
766 under Priestley-Taylor PET than Hamon PET and closely follows the process models in both cases.

767
768 The overall pattern of change in low flows (FLV) is very similar across all three DL models, with median
769 declines between -15% to -25% and little variability across sites (Figure 7c,d). The process models disagree
770 ~~significantly~~ on ~~the sign of~~ changes ~~to~~ ~~for~~ FLV, and ~~also~~ bound the changes predicted by the DL models.

771 HBV and HYMOD show mostly increases to FLV under warming and Priestley-Taylor PET, and a mix of
772 increases and decreases across sites for Hamon PET. SAC-SMA exhibits large declines in FLV under
773 warming and Hamon PET, and shows a median change that is similar to the DL models under Priestley-
774 Taylor PET. The percent changes in FLV across models tend to be large because the absolute magnitude
775 of FLV is small, and so small changes in millimeters of flow lead to large percent changes. [This can be
776 seen in sample daily hydrographs for two sites \(see Figure S5\), where visually the changes in low flows are
777 difficult to discern because they are all near zero for all models, but the change in the FLV statistic varies
778 significantly across the six models and two sites \(-56% to +40%\).](#)

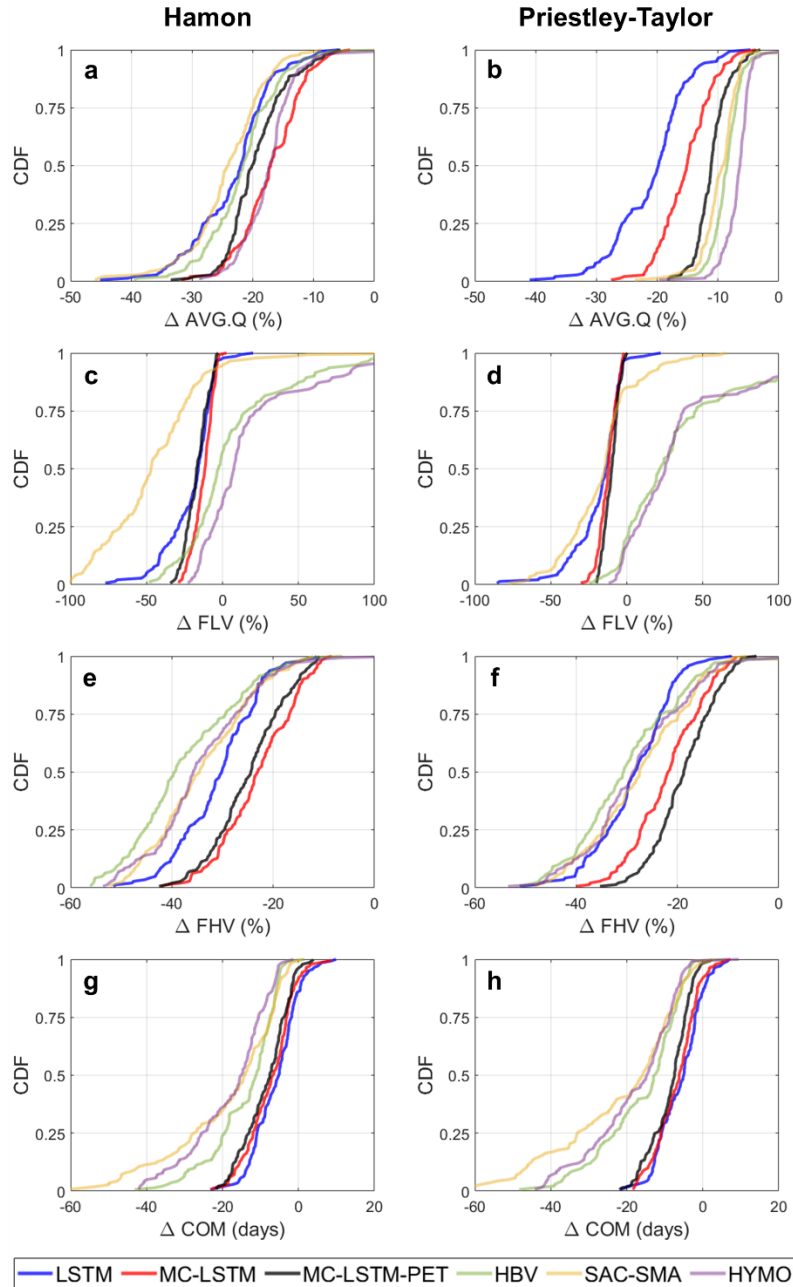
779
780 The differences between process-based and DL simulated changes for high flows (FHV; Figure 7e,f) and
781 [streamflow-seasonal](#) timing (COM; Figure 7g,h) are relatively consistent, with the process models
782 exhibiting [larger-more substantial](#) declines in high flows and earlier shifts in [streamflow-seasonal](#) timing
783 compared to the DL models. The choice of PET method has an [moderate](#)-impact on process-model based
784 changes in FHV, with larger declines under Hamon PET. A similar signal is also seen for the MC-LSTM-
785 PET but not the MC-LSTM or LSTM, although the LSTM predicts changes in FHV closest to the process
786 models.

787
788 For COM, the process models show a wide range of variability in projected change across sites, from no
789 change to 60 days earlier. For the DL models the range of change is much narrower, and the median change
790 in COM is [almost-approximately](#) a week less than the median change across the process models. [The earlier
791 shift in COM across all models is consistent with anticipated changes to snow accumulation and melt
792 dynamics under warming, with more water entering the stream during the winter and early spring as
793 precipitation shifts more towards rainfall and existing-snowpack melts off earlier in the year \(Byun and
794 Hamlet, 2018; Mote et al., 2018; Kayastha et al., 2022REFERENCES\). However, this effect is seen more
795 dramatically in the process models, as evidenced by more prominent changes to their daily and monthly
796 hydrographs under warming during the winter and early spring as compared to the DL models \(see see](#)

797 [Figures S5 and S6X](#)). The method of PET estimation has relatively little impact on both process model and
798 DL based estimates of change in COM.

799
800
801
802 [We note that the results above do not change even when considering the parametric uncertainty in the](#)
803 [process models, although for some metrics \(FLV\), uncertainty in process model estimated changes due to](#)
804 [parametric uncertainty is large \(see Figure S7\)](#). We [also](#) note that if the static watershed properties
805 (pet_mean, aridity, t_mean, frac_snow; see Table 1) are ~~also~~ changed to reflect warmer temperatures and
806 higher PET, all three DL models exhibit unrealistic water gains for between 15%-40% of locations
807 depending on the model and PET method, with the most water gains occurring under the LSTM (Figure
808 [S84](#)). These results suggest that changing the static watershed properties associated with long-term climate
809 characteristics can degrade the quality of the [projectionsestimated responses](#), at least when the ~~climate~~
810 [temperature changes shifts](#) are large and the range of average temperature and PET in the training set is
811 limited. ~~We also note that the results in Figure 7 are largely unchanged if based on projections for testing~~
812 ~~sites in the testing period (Figure S5)~~.

813



814

815 **Figure 7.** The distribution of change in (a,b) [long term mean daily flow \(AVG.Q\)](#), (c,d) [low flows \(FLV\)](#),
 816 (e,f) [high flows \(FHV\)](#), and (g,h) [seasonal streamflow timing \(COM\)](#) across the 141 training sites and all
 817 models under a scenario of 4°C warming using (a,c,e,g) Hamon PET and (b,d,f,h) Priestley-Taylor PET.

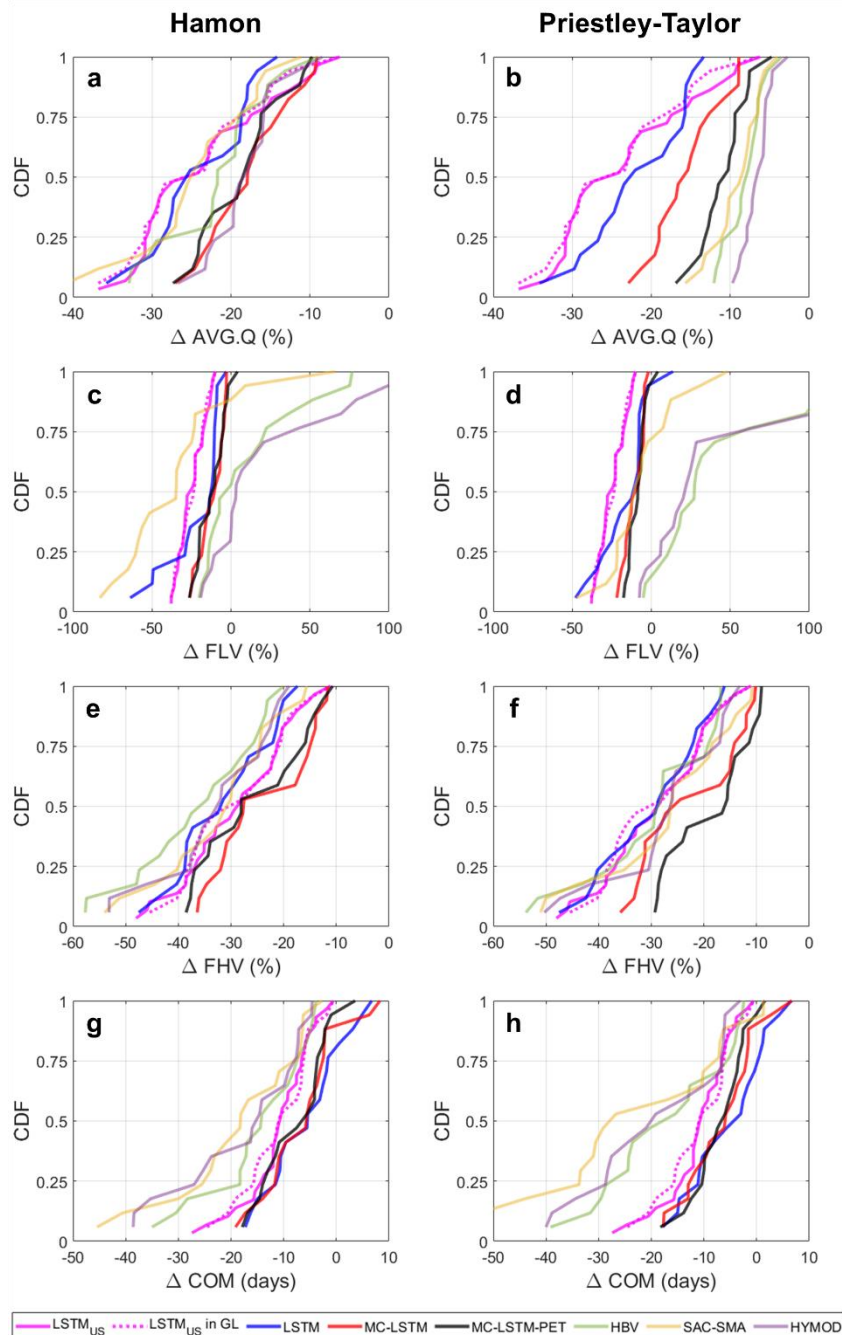
818 For the [DL-deep learning](#) models, changes were only made to the dynamic inputs (i.e., no changes to
 819 static inputs).
 820

821 One reason why the Great Lakes LSTM exhibits excessive [hydrologic-water](#) losses under warming could
 822 be that the model was trained using sites that are confined to a limited range of temperature and PET values

823 found in the Great Lakes basin (spanning approximately 40.5°-50°N), and so is ill-suited to extrapolate
824 hydrologic response under warming conditions that extend beyond ~~this range~~ temperature and PET range.
825 To evaluate this hypothesis, we examine changes to AVG.Q, FLV, FHV, and COM under 4°C warming at
826 the 29 CAMELS watersheds within the Great Lakes basin using the National LSTM (Figure 8). For
827 comparison, we also examine similar changes under all six Great Lakes DL and process models at 17 of
828 those 29 CAMELS basins that were used in the training and testing sets for the Great Lakes models. ~~;~~
829 ~~and~~ We also separate out highlight the National LSTM ~~projections-predictions~~ for those 17 sites. Note that
830 in Figure 8, the National LSTM ~~projections-predictions~~ do not differ between Hamon and Priestley Taylor
831 PET, because PET is not an input to that model.

832
833 The National LSTM was trained to watersheds across the CONUS (spanning approximately 26°-49°N),
834 and so was exposed to watersheds with much warmer conditions and higher PET during training. However,
835 we find that the National LSTM still ~~projects-predicts~~ very large declines in AVG.Q. For the 29 CAMELS
836 watersheds in the Great Lakes basin, the median decline in AVG.Q under the National LSTM is
837 approximately 25%, which is ~~only 0-6~~ moderately XX% larger than the median ~~projections-predictions~~ of
838 loss under the process models using Hamon PET ~~and but much XX~~ 16-19% larger than the process model
839 losses under Priestley-Taylor PET (Figure 8a,b). We also see larger declines in FLV under the National
840 LSTM as compared to the other Great Lakes DL models (Figure 8c,d). The National LSTM ~~projects-predicts~~
841 changes in FHV (Figure 8e,f) and COM (Figure 8g,h) that are relatively similar to the process models. ~~;~~ ~~and~~
842 ~~f~~ For COM, the projections-predictions of change are closer still smaller than to the process models but closer
843 to the process models than ~~for~~ any Great Lakes DL model. ~~;~~ suggesting that the National LSTM predicts
844 shifting snow accumulation and melt dynamics more consistently with the process models than regionally
845 fit DL models. In addition, the hydrologic ~~projections-predictions~~ are stable under the National LSTM
846 regardless of whether only dynamic inputs or both dynamic and static inputs are changed under warming
847 (see Figure S96), in contrast to the Great Lakes DL models. Therefore, the use of more watersheds in

848 training than span a more diverse set of climate conditions likely benefit the model when inputs are shifted
 849 significantly to reflect new climate conditions. However, as shown in Figure 8a,b, this benefit does not
 850 mitigate the tendency for the National LSTM to overestimate water loss under warming.



851

852 **Figure 8.** The distribution of change in (a,b) long term mean daily flow (AVG.Q), (c,d) low flows (FLV),
 853 (e,f) high flows (FHV), and (g,h) seasonal streamflow timing (COM) across 29 CAMELS sites within the
 854 Great Lakes basin under the National LSTM (solid pink), as well as for 17 of those 29 sites from the
 855 Great Lakes DL-deep learning and process models, under a scenario of 4°C warming. Results from the

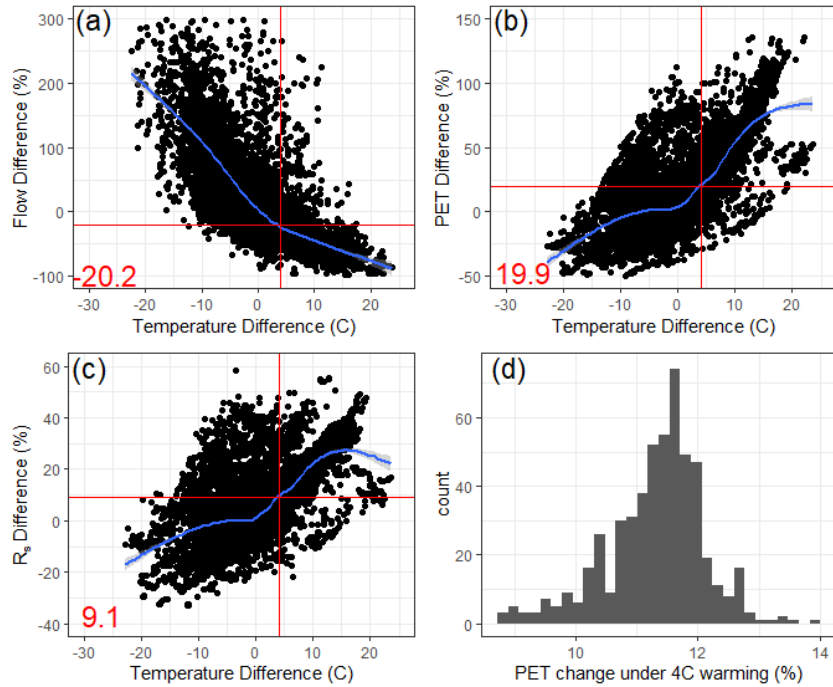
856 National LSTM for those 17 sites are also highlighted (dashed pink). For the Great Lakes models only,
857 results differ when using (a,c,e,f) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the National
858 LSTM, changes were made only to the dynamic inputs.
859

860 To better understand why the National LSTM predicts large water losses under warming, it is instructive
861 to examine how [average-long-term mean](#) streamflow, (Priestly-Taylor estimated) PET, and R_s vary across
862 all 531 CAMELS watersheds of different average temperatures, and compare this variability to [projected](#)
863 [predicted](#) changes in PET at each site under warming. Specifically, we ~~compare~~ [calculate](#) the difference in
864 long-term (1980-2014) [average-mean](#) streamflow (Figure 9a), PET (Figure 9b), and R_s (Figure 9c) across
865 all pairs of basins in the CAMELS dataset with average long-term precipitation within 1% of each other
866 [\(i.e., we only examine pairs of basins with very similar long-term mean precipitation\)](#). Then, for each basin
867 ~~pair, we~~ [and](#)-plot these differences [in long-term mean streamflow, PET, and \$R_s\$](#) against the differences in
868 [long-term](#) average temperature ~~across-for each-that~~ pair. The results show that the difference in [average](#)
869 [long-term mean](#) streamflow across watersheds with similar precipitation becomes negative when the
870 difference in temperature is positive (i.e., warmer watersheds have less flow on average), and that when the
871 difference in average temperature reaches 4°C, flows differ by about 20% on average (Figure 9a). This is
872 very similar to the ~~projected-predicted~~ median decline in [average-long-term mean](#) streamflow seen for the
873 National LSTM in Figure 8. We also note that average PET increases by approximately 20% between
874 watersheds that differ in average temperature by 4°C (Figure 9b). However, higher PET in warmer
875 watersheds is related both to the direct effect of temperature on vapor pressure deficit, as well as to the fact
876 that higher incoming solar radiation co-occurs in warmer watersheds (R_s is approximately 9% higher across
877 watershed pairs that differ by 4°C; Figure 9c). Using the Priestley-Taylor method, we estimate that average
878 PET would only increase by between 9-14% (median of 11.5%) if temperatures warm by 4°C and R_s is held
879 at historic values, while R_n is increased slightly due to declines in net outgoing longwave radiation with
880 warming (Figure 9d). However, the National LSTM appears to convolute the effects of temperature and R_s
881 and cannot separate out their effects on ~~ET-based~~ [evaporative](#) water loss, leading to larger [projected](#)
882 [predicted](#) streamflow losses under 4°C warming than changes in PET would warrant. This is possibly

883 because of the very strong correlation between at-site daily temperature and R_s historically (median
884 correlation of 0.85 across all CAMELS watersheds).

885

886



887

888 **Figure 9.** The percent difference in long-term (1980–2014) average-mean (a) streamflow, (b) Priestley-
889 Taylor based PET, and (c) downward shortwave radiation (R_s) for all pairs of CAMELS basins with
890 average precipitation within 1% of each other, plotted against differences in average temperature for each
891 pair. A loess smooth is provided for each scatter (blue), along with the changes in variable estimated at a
892 4°C temperature difference between pairs of sites (red). (d) The projected change in Priestley-Taylor
893 based PET (as a percentage) for each CAMELS basin under 4°C warming, assuming no change in R_s .
894

895 5. Discussion and Conclusion

896 In this study, we contribute a sensitivity analysis that evaluates the physical plausibility of future
897 streamflow projections-responses under climate-changewarming using DL rainfall-runoff models. The basis
898 for this evaluation is anchored to the assumption that differences in estimated streamflow projections
899 responses should emerge under very different projections-scenarios of future-PET under warming, and that
900 realistic projections-predictions of future-PET and water loss under warming tend to be much lower than
901 those estimated by temperature-based PET methods. Accordingly, we assume that physically plausible

902 ~~future~~ streamflow ~~projections~~ ~~predictions~~ should be able to respond to lower energy-budget based PET
903 projections under warming and, all else equal, ~~project~~ ~~estimate~~ smaller streamflow losses.

904
905 The results of this study show that a standard LSTM is not able to predict physically realistic differences in
906 streamflow response across substantially different ~~projections~~ ~~estimates~~ of ~~future~~ PET under warming. This
907 discrepancy ~~in future projections~~ emerged despite the fact that the standard LSTM was a far better model
908 for streamflow estimation in ungauged basins compared to three process-based models under historic
909 climate conditions. In addition, the National LSTM trained to a much larger set of watersheds (531 basins
910 across 23° of latitude) using temperature, vapor pressure, and R_s directly (rather than PET) also estimated
911 water loss under warming that far exceeded the losses estimated with process models forced with energy
912 budget-based PET. Since water losses estimated using energy budget-based PET are generally considered
913 more realistic (Lofgren et al., 2011; Shaw and Riha, 2011; Lofgren and Rouhana, 2016; Milly and Dunne,
914 2017; Lemaitre-Basset et al. 2022), this result casts doubt over the physical plausibility of the LSTM
915 ~~projection~~ ~~predictions~~.

916
917 Results from this work also suggest that PIML-based DL models can capture physically plausible
918 streamflow responses under ~~climate change~~ ~~warming~~ while still maintaining superior prediction skill
919 compared to process models, at least in some cases. In particular, a mass conserving LSTM that also
920 respected the limits of water loss due to ~~ET~~ ~~evapotranspiration~~ (the MC-LSTM-PET) was able to ~~project~~
921 ~~predict~~ changes in ~~average~~ ~~long-term~~ ~~mean~~ streamflow that much more closely aligned with process-model
922 based estimates, while also providing competitive out-of-sample performance across all models considered
923 (including the other DL models). A more conventional MC-LSTM that did not limit water losses by PET
924 was less consistent with process-based estimates of change in ~~average~~ ~~long-term~~ ~~mean~~ streamflow. These
925 results highlight the potential for PIML-based DL models to help achieve similar performance
926 improvements over process-based models as documented in recent work on DL rainfall-runoff models

927 (Kratzert et al., 2019a,b; Feng et al., 2020; Nearing et al., 2021) while also producing projections under
928 climate change that are more consistent with theory than non-PIML DL models.

929
930 An interesting result from this study was the disagreement in the change in high flows and [streamflow](#)
931 [seasonal streamflow](#) timing between all Great Lakes DL models and process models, the latter which
932 estimated greater reductions in high flows and larger shifts of water towards earlier in the year. [Projections](#)
933 [Predictions](#) from the Great Lakes DL models were also unstable if static climate properties of each
934 watershed were changed under warming. In contrast, the National LSTM was more stable if static properties
935 were changed, and it predicted changes to high flows and [streamflow-seasonal](#) timing that were more like
936 the process models than [projections-predictions](#) from the Great Lakes DL models. [The results for COM in](#)
937 [particular suggest that the National LSTM is may be more consistent with the process models in terms of](#)
938 [its representation of warming effects on snow accumulation and melt processes and the resulting shifts in](#)
939 [the seasonal hydrograph, although differences with the process model predictions were still notable. Still,](#)
940 [these results are consistent with past work showing that large-sample LSTMs can learn to represent snow](#)
941 [processes internally from meteorological and streamflow data \(Lees et al., 2022\).](#) While its challenging to
942 know which set of [projections-predictions](#) are correct for these streamflow properties, these results overall
943 favor [projections-predictions](#) from the National LSTM [over the regional LSTMs](#) and highlight the benefits
944 of DL rainfall-runoff models trained to a larger set of diverse watersheds for climate change analysis.

945
946 [To properly interpret the results of this work, there are several limitations of this study that require](#)
947 [discussion. First there were differences in the inputs and data sources between the National LSTM and all](#)
948 [other Great Lakes models, including the source of meteorological data and the lack of PET as an input into](#)
949 [the National LSTM. While this latter discrepancy might be less impactful \(i.e., the National LSTM was](#)
950 [provided meteorological inputs that together completely determine Hamon and Priestley-Taylor PET\), the](#)
951 [difference in meteorological data across the two sets of models is a substantial source of uncertainty and](#)
952 [could lead to non-trivial differences in hydrologic response estimation, complicating a direct comparison](#)

953 [of the National LSTM to the other models. Future work for the Great Lakes Intercomparison Project should](#)
954 [consider developing consistent datasets with other \(and larger\) benchmark datasets like CAMELS to](#)
955 [address this issue.](#)

956 ~~The MC-LSTM-PET model proposed in this work represents one (relatively simple) PIML-based~~
957 ~~architectural change to an existing DL model in the hydrologic literature that can help better capture~~
958 ~~physical constraints on water loss from hydrologic systems. However, other possibilities exist. For example,~~
959 ~~the hard constraint in the MC-LSTM-PET could instead be imposed as a soft constraint through adjustments~~
960 ~~to the loss function, where water losses in the trash cell that exceed PET are penalized. The MC-LSTM-~~
961 ~~PET model could also be adjusted further to allow additional water losses in the trash cell related to human~~
962 ~~water extractions from the watershed or other terminal sinks. A different approach would be to use learnable,~~
963 ~~differentiable, process-based models with embedded neural networks (Jiang et al., 2020; Feng et al., 2022;~~
964 ~~Feng et al., 2023), which can achieve similar performance to LSTMs but can also represent and output~~
965 ~~different internal hydrologic fluxes. Further work is needed to evaluate the benefits and drawbacks of these~~
966 ~~different PIML-based approaches, preferably on large benchmarking datasets such as CAMELS.~~

967
968
969 [OneAnother](#) important limitation [of this study](#) is how we constructed the [climate change warming](#) scenarios,
970 with 4°C warming [and shifts to PET but but](#) no changes [to net incoming shortwave radiation and slight](#)
971 [decreases in net outgoing longwave radiation with warming \(i.e., slight increases in \$R_n\$ \) to other](#)
972 [meteorological variables \(net incoming shortwave radiation, precipitation, humidity, air pressure, wind](#)
973 [speeds\). These scenarios and associated sensitivity analyses were constructed in the style of other](#)
974 [metamorphic tests for hydrologic models \(Yang and Chui, 2021; Razavi, 2021; Reichert et al., 2023\), where](#)
975 [we define input changes with expected responses and test whether model behavior is consistent with these](#)
976 [expectations. However, for DL and other machine learning-\(ML\) models, the results of such sensitivity](#)
977 [analyses may be unreliable because of distributional shifts between the training and testing data and poor](#)
978 [out-of-distribution generalization \(see Shen et al., 2021, Wang et al., 2023, and references within\). When](#)

979 trained, conventional machine learning ML-models try to leverage all of the correlations within the training
980 set to minimize training errors, which is effective in out-of-sample performance only if those same patterns
981 of correlation persistent into the testing data (Liu et al., 2021). In our experimental design, we impose a
982 distinct shift in the joint distribution of the inputs (i.e., a covariate shift) by increasing temperatures and
983 PET but leaving unchanged other meteorological inputs, thereby altering the correlation among inputs.
984 Therefore, one might expect some degradation in the DL model-based predictions of streamflow under
985 these scenarios.

986
987 ~~While outside the scope of the present study, we~~The challenge of out-of-distribution generalization and its
988 application to DL rainfall-runoff model testing under climate change highlights several important avenues
989 for future work. First, additional efforts are needed to evaluate the ~~argue more work is needed to further~~
990 ~~explore the physical plausibility of DL-based hydrologic projections under climate change~~with more
991 ~~standard~~ while ensuring that LSTMs, with greater attention paid to the joint distribution of all
992 ~~meteorologiemeteorological inputs used in future scenarios is realistic. For example, there are physical~~
993 relationships between changes in temperature and net radiation (Nordling et al., 2021), as well as
994 temperature, humidity, and extreme precipitation (Ali et al., 2018; Najibi et al., 2022), that should all be
995 preserved in future climate scenarios. The use of climate model output may be well suited for such tests,
996 although care is needed to avoid significant statistical bias correction and downscaling (i.e., post-processing)
997 of multiple climate fields that could cause shifts in the joint distribution across inputs (Maraun, 2016). High-
998 resolution convective-permitting models may be helpful in this regard, given their improved accuracy for
999 key climate fields like precipitation ((Kendon et al. 2017).

1000
1001
1002 ~~the model under historical and future climate conditions. We did not consider any changes in net incoming~~
1003 ~~shortwave radiation because there is significant uncertainty in this term at local scales and its relationship~~
1004 ~~to local temperature change. Projections of net incoming shortwave radiation are highly variable across~~

1005 space and can even differ in the direction of change, largely because of uncertainty in the representation of
1006 clouds in climate models, future projections of aerosols, and the representation of cloud-aerosol interactions
1007 (Chen, 2021; Coppola et al., 2021; Taranu et al., 2023). The relationship between local net radiation change
1008 and local temperature change further depends on horizontal energy transport from other regions (Nordling
1009 et al., 2021). In addition, the approximation we used for changes to net outgoing longwave radiation was
1010 not designed to resolve all land-atmosphere energy balance feedbacks with changing atmospheric
1011 composition under climate change. These uncertainties, along with uncertainties in energy budget based
1012 methods used to estimate PET (Greve et al. 2019; Liu et al., 2022), complicate future projections of
1013 atmospheric drying power under warming. Regardless, the main finding of this work remains, namely that
1014 DL models struggle to propagate different hypotheses of future PET scenarios into hydrologic projections
1015 unless explicitly directed to do so.

1016 There are also several emerging techniques in machine learning ML to address out-of-distribution
1017 generalization directly (Shen et al., 2021). One family set of promising methods for the challenge of DL
1018 hydrologic modeling under climate change is causal learning, defined broadly as methods that aimed to
1019 identifying input variables that have a causal relationship with the target variable and to leverage those
1020 inputs for prediction (Shen et al., 2021). PIML One approach for this is to approaches, such as the MC-
1021 LSTM-PET model proposed in this work, fall into this category (Vasudevan et al., 2021). Here, prior
1022 scientific knowledge on causal structures can be embedded into the DL model through tailored loss
1023 functions or, as in the case of the MC-LSTM-PET model, through ~~The MC LSTM PET model proposed~~
1024 ~~in this work represents one (relatively simple) PIML-based architectural adjustments or constraints (for~~
1025 ~~other examples outside of hydrology, see Lin et al., 2017; Ma et al., 2018)~~ change to an existing DL model
1026 ~~in the hydrologic literature that can help better capture physical constraints on water loss from hydrologic~~
1027 ~~systems. The MC-LSTM-PET model can be viewed as a specific, limited case of a broader class of However,~~
1028 ~~other possibilities exist. For example, the hard constraint in the MC LSTM PET could instead be imposed~~
1029 ~~as a soft constraint through adjustments to the loss function, where water losses in the trash cell that exceed~~
1030 ~~PET are penalized. The MC LSTM PET model could also be adjusted further to allow additional water~~

1031 ~~losses in the trash cell related to human water extractions from the watershed or other terminal sinks. A~~
1032 ~~different approach would be to use learnable, differentiable, process-based models with embedded neural~~
1033 ~~networks (also referred to as hybrid differentiable models; Jiang et al., 2020; Feng et al., 2022; Feng et al.,~~
1034 ~~2023a). These models use process model architectures as a backbone for model structure, which is then~~
1035 ~~enhanced through flexible, data-driven learning for a subset of processes. Recent work has shown that these~~
1036 ~~models, which can achieve similar performance to LSTMs but can also represent and output different~~
1037 ~~internal hydrologic fluxes (Feng et al., 2022; Feng et al., 2023a).~~

1038

1039 However, challenges can arise when imposing architectural constraints in PIML models. For example, the
1040 MC-LSTM-PET model makes the assumption that all water loss in the system is due to evapotranspiration,
1041 and therefore cannot exceed PET. However, other terminal sinks are possible, such as human water
1042 extractions and inter-basin transfers (Siddik et al. 2023) or water lost to aquifer recharge and inter-basin
1043 groundwater fluxes (Safeeq et al., 2021; Jasechko et al., 2021). It is difficult to know the magnitude of these
1044 alternative sinks given unknown systematic errors in other inputs (e.g., underestimation of precipitation
1045 from under-catch) that confound water balance closure analyses. Still, recent techniques and datasets to
1046 help quantify these sinks (Gordon et al., 2022; Siddik et al. 2023) provide an avenue to integrate them into
1047 the MC-LSTM-PET model constraints to improve generalizability. However, Yet as constraints are added
1048 to the model architecture (i.e., more assumptions are inherited from a process model backbone), the
1049 potential grows for inductive bias that negatively impacts generalizability. For instance, a recent evaluation
1050 of hybrid differentiable models showed that they underperformed relative to a standard LSTM due to
1051 structural deficiencies in cold regions, arid regions, and basins with considerable anthropogenic impacts
1052 (Feng et al., 2023b). Some of these challenges may be difficult to address because only differentiable
1053 process models can be considered in this hybrid framework, limiting the process model structures that could
1054 be adapted with this approach. Further
1055 Additional work is needed to evaluate the benefits and drawbacks of
1056 these different PIML-based approaches, preferably on large benchmarking datasets such as CAMELS or
CAVARAN (Kratzert et al., 2023).

1057
1058 Given some of the potential challenges above,
1059 other DL methods that advance causality while making fewer assumptions on watershed-scale process
1060 controls are also worth pursuing. For example, a series of techniques have emerged that embed the concept
1061 and constraints of directed acyclic graphs within deep neural networks in such a way that the architecture
1062 of the neural network is inferred from the data to encode causality among variables (see Luo et al., 2020
1063 and references within). That is, frameworks to optimize the architecture of the model can be designed not
1064 only to maximize out-of-sample predictive performance, but also to promote causality. Alternatively,
1065 domain-invariant learning attempts to promote the identification of features that are domain-specific versus
1066 domain invariant, by separating and labeling training data from different ‘domains’ or ‘environments’ (Ilse
1067 et al., 2021). In the case of DL rainfall-runoff models, this strategy could be implemented, for instance, by
1068 pairing observed climate and streamflow (one domain) with land surface model-based streamflow estimated
1069 using future projected climate model output (another domain), with the goal to learn invariant relationships
1070 between key climate inputs (e.g., net radiation or PET) and streamflow across the two domains. Here, there
1071 may be a benefit from including data from the land surface and climate models, where the correlation
1072 between temperature, net radiation, and PET may be weaker under projected climate change. These
1073 techniques offer an intriguing alternative for the next generation of DL hydrologic models that can
1074 generalize well under climate change, and should be the focus of further exploration. ~~identify inputs where~~
1075 ~~the conditional distribution of the target variable (streamflow) given that input is invariant across~~
1076 ~~heterogeneous datasets. A large focus on~~

1077
1078
1079
1080 Finally, we note that the results of this study do not entirely preclude the possibility that a standard LSTM,
1081 fit to a sufficiently large set of diverse watersheds, could ultimately learn more physically realistic
1082 projections under climate change. Our results with the National LSTM suggest that the signals between

1083 temperature change and R_s on water loss may be entangled, making it difficult for the model to estimate the
1084 individual effects of changes to one of those terms (temperature) on water loss. However, it is possible that
1085 the model would produce hydrologic projections that were more in line with theory if it was given 1) high
1086 quality data on all terms related to water loss; and 2) future projections of these terms that were co-
1087 developed in physically consistent ways (e.g., from physical climate models). The R_s used in the National
1088 LSTM was based on reanalysis and so may have had meaningful errors that drove the model to attribute
1089 more water loss to warmer temperatures, and the scenario of warming given to the National LSTM (4°C
1090 warming with no change in R_s) may violate the physical relationship between temperatures and R_s . While
1091 outside the scope of the present study, we argue more work is needed to further explore the physical
1092 plausibility of hydrologic projections with more standard LSTMs, with greater attention paid to the
1093 meteorologic inputs used in the model under historical and future climate conditions.

1095 **Acknowledgements**

1096 This research was supported by the U.S. National Science Foundation grant CBET-2144332.

1098 **Data Availability Statement**

1099 The code used for this project is available at <https://doi.org/10.5281/zenodo.8190287> at
1100 <https://doi.org/10.5281/zenodo.10027355>. All data used to -train and evaluate the
1101 models are available at https://www.hydrohub.org/mips_introduction.html#grip-gl.

1103 **References**

1104 [Ali, H., Fowler, H. J., & Mishra, V. \(2018\). Global observational evidence of strong linkage between dew
1105 point temperature and precipitation extremes. Geophysical Research Letters, 45, 12320–
1106 12330. <https://doi.org/10.1029/2018gl080557>](#)

1107
1108 Allen, R.G., Pereira, L.S., Raes, D., et al. (1998) Crop Evapotranspiration-Guidelines for Computing
1109 Crop Water Requirements-FAO Irrigation and Drainage Paper 56. FAO, Rome, 300(9): D05109.

1111 Anderson, E. A. (1976). A point energy and mass balance model of a snow cover (NOAA Technical
1112 Report NWS 19). Silver Spring, MD: National Oceanic and Atmosphere Administration.
1113

1114 [Bastola S., Murphy C., Sweeney J. \(2011\). The role of hydrological modelling uncertainties in climate
1115 change impact assessments of Irish river catchments. *Adv Water Resour.*, 34, 562–76.](#)
1116

1117 [Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Miralles, D. G., McVicar, T. R., Schellekens, J.,
1118 and Bruijnzeel, L. A. \(2016\), Global-scale regionalization of hydrologic model parameters, *Water Resour.
1119 Res.*, 52, 3599–3622, doi:10.1002/2015WR018247.](#)

1120 [Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Dutra, E., Fink, G., Orth, R., and Schellekens, J.: Global
1121 evaluation of runoff from 10 state-of-the-art hydrological models \(2017\), *Hydrol. Earth Syst. Sci.*, 21,
1122 2881–2903, <https://doi.org/10.5194/hess-21-2881-2017>.](#)

1123 Bergström, S. & Forsman, A. (1973) Development of a conceptual deterministic rainfall-runoff model.
1124 *Nordic Hydrol.* 4, 147–170.
1125

1126 Beven, K. (2023). Benchmarking hydrological models for an uncertain future. *Hydrological
1127 Processes*, 37(5), e14882. <https://doi.org/10.1002/hyp.14882>
1128

1129 Boyle, D. P. (2001). Multicriteria calibration of hydrologic models, (Doctoral dissertation). Retrieved from
1130 UA Campus Repository (<http://hdl.handle.net/10150/290657>), Tucson, AZ: The University of Arizona.
1131

1132 [Breuer, L., Huisman, J. A., Willems, P., Bormann, H., Bronstert, A., Croke, B. F. W., Frede, H.-G., Gräff,
1133 T., Hubrechts, L., Jakeman, A. J., Kite, G., Lanini, J., Leavesley, G., Lettenmaier, D. P., Lindström, G.,
1134 Seibert, J., Sivapalan, M., and Viney, N. R.: Assessing the impact of land use change on hydrology by
1135 ensemble modeling \(LUCHEM\). I: Model intercomparison with current land use, *Adv. Water Resour.*,
1136 32, 129–146, <https://doi.org/10.1016/j.advwatres.2008.10.003>, 2009.](#)
1137

1138 Burnash, R. J. (1995). The NWS river forecast system - catchment modeling. In Singh, V. (Ed.), *Computer
1139 Models of Watershed Hydrology* (pp. 311-366). Littleton, CO: Water Resources Publication.
1140

1141

1142 [Byun, K. and Hamlet, A.F. \(2018\), Projected changes in future climate over the Midwest and Great Lakes
1143 region using downscaled CMIP5 ensembles. *Int. J. Climatol.*, 38: e531-
1144 e553. <https://doi.org/10.1002/joc.5388>](#)
1145

1146 Campbell, M., Cooper, M. J. P., Friedman, K., & Anderson, W. P. (2015). The economy as a driver of
1147 change in the Great Lakes - St. Lawrence basin. *Journal of Great Lakes Research*, 41, 69–83.
1148

1149 [Cayan, D. R., Kammerdiener, S. A., Dettinger, M. D., Caprio, J. M., & Peterson, D. H. \(2001\). Changes
1150 in the Onset of Spring in the Western United States, *Bulletin of the American Meteorological
1151 Society*, 82\(3\), 399–416. \[https://doi.org/10.1175/1520-0477\\(2001\\)082<0399:CITOOS>2.3.CO;2\]\(https://doi.org/10.1175/1520-0477\(2001\)082<0399:CITOOS>2.3.CO;2\)
1152](#)

1153

1154 [Clark, M. P., Bierkens, M. F. P., Samaniego, L., Woods, R. A., Uijlenhoet, R., Bennett, K. E., Pauwels,
1155 V. R. N., Cai, X., Wood, A. W., and Peters-Lidard, C. D. \(2017\). The evolution of process-based
1156 hydrologic models: historical challenges and the collective quest for physical realism, *Hydrol. Earth Syst.
1157 Sci.*, 21, 3427–3440, <https://doi.org/10.5194/hess-21-3427-2017>.](#)

1158 [Clark, M.P., Wilby, R.L., Gutmann, E.D. et al. Characterizing Uncertainty of the Hydrologic Impacts of](#)
1159 [Climate Change. *Curr Clim Change Rep* 2, 55–64 \(2016\). <https://doi.org/10.1007/s40641-016-0034-x>](#)
1160 [Chen, L. Uncertainties in solar radiation assessment in the United States using climate models. *Clim*
1161 \[Dyn\]\(#\) 56, 665–678 \(2021\). <https://doi.org/10.1007/s00382-020-05498-7>](#)
1162
1163
1164 [Coxon, G., Addor, N., Bloomfield, J. P., Freer, J., Fry, M., Hannaford, J., Howden, N. J. K., Lane, R.,](#)
1165 [Lewis, M., Robinson, E. L., Wagener, T., and Woods, R. \(2020\). CAMELS-GB: hydrometeorological](#)
1166 [time-series and landscape attributes for 671 catchments in Great Britain, *Earth Syst. Sci. Data*, 12, 2459–](#)
1167 [2483, <https://doi.org/10.5194/essd-12-2459-2020>.](#)
1168
1169 [Coppola, E., Nogherotto, R., Ciarlò, J. M., Giorgi, F., van Meijgaard, E., Kadyrov, N., et al.](#)
1170 [\(2021\). Assessment of the European Climate Projections as Simulated by the Large EURO-CORDEX](#)
1171 [Regional and Global Climate Model Ensemble. *Journal of Geophysical Research: Atmospheres*, 126,](#)
1172 [e2019JD032356. <https://doi.org/10.1029/2019JD032356>](#)
1173
1174 Demargne, J. et al. (2014). The Science of NOAA's Operational Hydrologic Ensemble Forecast
1175 Service. *Bull. Amer. Meteor. Soc.*, 95, 79–98, <https://doi.org/10.1175/BAMS-D-12-00081.1>.
1176
1177 [Fan, Y. \(2019\). Are catchments leaky? *WIREs Water*, 6\(6\). <https://doi.org/10.1002/wat2.1386>](#)
1178
1179 Feng, D., Fang, K., & Shen, C. (2020). Enhancing streamflow forecast and extracting insights using long-
1180 short term memory networks with data integration at continental scales. *Water Resources Research*, 56,
1181 e2019WR026793. <https://doi.org/10.1029/2019WR026793>
1182
1183 Feng, D., Liu, J., Lawson, K., & Shen, C. (2022). Differentiable, learnable, regionalized process-based
1184 models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy. *Water*
1185 *Resources Research*, 58, e2022WR032404. <https://doi.org/10.1029/2022WR032404>
1186
1187 Feng, D., Beck, H., Lawson, K., and Shen, C. (2023a). The suitability of differentiable, physics-informed
1188 machine learning hydrologic models for ungauged regions and climate change impact assessment,
1189 *Hydrol. Earth Syst. Sci.*, 27, 2357–2373, <https://doi.org/10.5194/hess-27-2357-2023>.
1190
1191 [Feng, D., Beck, H., de Bruijn, J., Sahu, R. K., Satoh, Y., Wada, Y., Liu, J., Pan, M., Lawson, K., and](#)
1192 [Shen, C. \(2023b\). Deep Dive into Global Hydrologic Simulations: Harnessing the Power of Deep](#)
1193 [Learning and Physics-informed Differentiable Models \(\$\delta\$ HBV-globe1.0-hydroDL\), *Geosci. Model Dev.*](#)
1194 [Discuss. \[preprint\], <https://doi.org/10.5194/gmd-2023-190>, in review.](#)
1195
1196 Frame, J.M., Kratzert, F., Gupta, H.V., Ullrich, P., & Nearing, G.S. (2022). On Strictly enforced mass
1197 conservation constraints for modeling the Rainfall-Runoff process. *Hydrological Processes*, 37, e14847,
1198 <https://doi.org/10.1002/hyp.14847>.
1199
1200 Frame, J.M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., et al. (2021b). Deep learning
1201 rainfall-runoff predictions of extreme events. *Hydrology and Earth System Sciences*, 26, 3377-
1202 3392, <https://doi.org/10.5194/hess-26-3377-2022>.
1203
1204 Frame, J.M., Kratzert, F., Raney II, A., Rahman, M., Salas, F.R., & Nearing, G.S. (2021a). Post-
1205 processing the National Water Model with Long Short-Term Memory networks for streamflow
1206 predictions and diagnostics. *Journal of the American Water Resources Association*, 1-12.
1207 <https://doi.org/10.1111/1752-1688.12964>
1208

1209 Fry, L. M., Hunter, T. S., Phanikumar, M. S., Fortin, V., and Gronewold, A. D. (2013), Identifying
1210 streamgauge networks for maximizing the effectiveness of regional water balance modeling, *Water Resour.*
1211 *Res.*, 49, 2689– 2700, doi:10.1002/wrcr.20233.
1212

1213 Gasset, N., Fortin, V., Dimitrijevic, M., Carrera, M., Bilodeau, B., Muncaster, R., Gaborit, É., Roy, G.,
1214 Pentcheva, N., Bulat, M., Wang, X., Pavlovic, R., Lespinas, F., Khedhaouiria, D., and Mai, J.: A 10 km
1215 North American precipitation and land-surface reanalysis based on the GEM atmospheric model, *Hydrol.*
1216 *Earth Syst. Sci.*, 25, 4917–4945, <https://doi.org/10.5194/hess-25-4917-2021>, 2021.
1217

1218 Gauch, M., Kratzert, F., Klotz, D., Nearing, G., Lin, J., & Hochreiter, S. (2021a). Rainfall-runoff
1219 prediction at multiple timescales with a single Long Short-Term Memory network. *Hydrology and Earth*
1220 *System Sciences*, 25, 2045-2062. <https://doi.org/10.5194/hess-25-2045-2021>
1221

1222 Gauch, M., Mai, J., & Lin, J. (2021b). The proper care and feeding of CAMELS: How limited training
1223 data affects streamflow prediction. *Environmental Modelling and Software*, 135, 104926.
1224 <https://doi.org/10.1016/j.envsoft.2020.104926>
1225

1226 ~~Greve, P., Roderick, M.L., Ukkola, A.M., and Wada, Y. (2019), The aridity index under global warming,~~
1227 ~~*Environmental Research Letters*, 14, 124006, <https://doi.org/10.1088/1748-9326/ab5046>.~~
1228

1229 ~~Gordon, B.L., Brooks, P.D., Krogh, S.A., Boisrime, G.F.S., Carrol, R.W.H., McNamara, J.P., & Harpold,~~
1230 ~~A.A. (2022), Why does snowmelt driven streamflow response to warming vary? A data driven review~~
1231 ~~and predictive framework, *Environmental Research Letters*, 15 (5), 053004. [https://doi.org/10.1088/1748-](https://doi.org/10.1088/1748-9326/ae64b4)~~
1232 ~~[9326/ae64b4](https://doi.org/10.1088/1748-9326/ae64b4)~~
1233 ~~Gordon, B. L., Crow, W. T., Konings, A. G., Dralle, D. N., & Harpold, A. A. (2022). Can we use the water~~
1234 ~~budget to infer upland catchment behavior? The role of data set error estimation and interbasin~~
1235 ~~groundwater flow. *Water Resources Research*, 58, e2021WR030966. [https://](https://doi.org/10.1029/2021WR030966)~~
1236 ~~doi.org/10.1029/2021WR030966~~
1237

1238 ~~Greve, P., Roderick, M.L., Ukkola, A.M., and Wada, Y. (2019), The aridity index under global warming,~~
1239 ~~*Environmental Research Letters*, 14, 124006, <https://doi.org/10.1088/1748-9326/ab5046>.~~
1240

1241 Gronewold, A. D., and Rood, R. B. (2019). Recent water level changes across Earth’s largest lake system
1242 and implications for future variability. *Journal of Great Lakes Research*, 45(1), 1–3.
1243

1244 Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F. (2009). Decom- position of the mean squared
1245 error and NSE performance criteria: Implications for improving hydrological modelling, *J. Hydrol.*, 377,
1246 80–91.
1247

1248 Hamon, W. R. (1963). Estimating Potential Evapotranspiration, *T. Am. Soc. Civ. Eng.*, 128, 324–
1249 338, <https://doi.org/10.1061/TACEAT.0008673>.
1250

1251 Hansen, C., Shafiei Shiva, J., McDonald, S., and Nabors, A. (2019). Assessing Retrospective National
1252 Water Model Streamflow with Respect to Droughts and Low Flows in the Colorado River Basin. *Journal*
1253 *of the American Water Resources Association* 964– 975. <https://doi.org/10.1111/1752-1688.12784>.
1254

1255 ~~Hargreaves, G.H., and Samani, Z.A. (1985). Reference crop evapotranspiration from~~
1256 ~~temperature. *Applied Engineering in Agriculture* 1: 96–99.~~
1257

1258 Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-
1259 1780. <https://doi.org/10.1162/neco.1997.9.8.1735>

1260
1261 Hoedt, P.J., F. Kratzert, D. Klotz, C. Halmich, M. Holzleitner, G. Nearing, et al. (2021). MC-LSTM:
1262 Mass-Conserving LSTM. *arXiv e-prints*, arXiv:2101.05186. Retrieved from
1263 <https://arxiv.org/abs/2101.05186>
1264
1265 [Höge, M., Scheidegger, A., Baity-Jesi, M., Albert, C., and Fenicia, F. \(2022\). Improving hydrologic](#)
1266 [models for predictions and process understanding using neural ODEs, *Hydrol. Earth Syst. Sci.*, 26, 5085–](#)
1267 [5102, <https://doi.org/10.5194/hess-26-5085-2022>.](#)
1268
1269 Hrachowitz, M. et al. (2013). A decade of Predictions in Ungauged Basins (PUB)—a
1270 review, *Hydrological Sciences Journal*, 58:6, 1198-1255, DOI: 10.1080/02626667.2013.803183
1271
1272
1273
1274 [Ilse, M., Tomczak, J.M., and Forré, P. \(2021\). Selecting Data Augmentation for Simulating Interventions.](#)
1275 [Proceedings of the 38th International Conference on Machine Learning, PMLR 139:4555-4562.](#)
1276
1277 [Jasechko, S., Seybold, H., Perrone, D. et al. Widespread potential loss of streamflow into underlying](#)
1278 [aquifers across the USA. *Nature* 591, 391–395 \(2021\). <https://doi.org/10.1038/s41586-021-03311-x>](#)
1279
1280 Jiang, S., Zheng, Y., & Solomatine, D. (2020). Improving AI system awareness of geoscience knowledge:
1281 Symbiotic integration of physical approaches and deep learning. *Geophysical Research Letters*, 46,
1282 e2020GL088229. <https://doi.org/10.1029/2020GL088229>
1283
1284 ~~[Kapnick, S., & Hall, A. \(2010\). Observed Climate–Snowpack Relationships in California and their](#)~~
1285 ~~[Implications for the Future, *Journal of Climate*, 23\(13\), 3446–](#)~~
1286 ~~[3456. <https://doi.org/10.1175/2010JCLI2903.1>](#)~~
1287
1288 Karpantne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., et al. (2017).
1289 Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Transactions on*
1290 *Knowledge and Data Engineering*, 29(10), 2318-2331. <https://doi.org/10.1109/TKDE.2017.2720168>
1291
1292 Kayastha, M.B., Ye, X., Huang, C., and Xue, P. (2022), Future rise of the Great Lakes water levels under
1293 climate change, *Journal of Hydrology*, 612 (Part B), 128205,
1294 <https://doi.org/10.1016/j.jhydrol.2022.128205>.
1295
1296
1297
1298 [Kendon, Elizabeth J., Nikolina Ban, Nigel M. Roberts, Hayley J. Fowler, Malcolm J. Roberts, Steven C.](#)
1299 [Chan, Jason P. Evans, Giorgia Fosser, and Jonathan M. Wilkinson. \(2017\). Do Convection-Permitting](#)
1300 [Regional Climate Models Improve Projections of Future Precipitation Change? *Bulletin of the American*](#)
1301 [Meteorological Society 98 \(1\): 79–93. <https://doi.org/10.1175/BAMS-D-15-0004.1>.](#)
1302
1303 Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *arXiv e-prints*,
1304 arXiv:1412.6980. Retrieved from <https://arxiv.org/abs/1412.6980>
1305
1306 [Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S.,](#)
1307 [and Nearing, G. \(2022\). Uncertainty estimation with deep learning for rainfall–runoff modeling, *Hydrol.*](#)
1308 [Earth Syst. Sci.](#), 26, 1673–1693, <https://doi.org/10.5194/hess-26-1673-2022>.

1309 Konapala, G., Kao, S. C., Painter, S., & Lu, D. (2020). Machine learning assisted hybrid models can
1310 improve streamflow simulation in diverse catchments across the conterminous US. *Environmental*
1311 *Research Letters*, 15(10), 104022. <https://doi.org/10.1088/1748-9326/aba927>
1312

1313 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019a).
1314 Toward improved predictions in ungauged basins: Exploiting the power of machine learning. *Water*
1315 *Resources Research*, 55, 11,344–11,354. <https://doi.org/10.1029/2019WR026065>
1316

1317 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. S. (2019b). Towards
1318 learning universal, regional, and local hydrological behaviors via machine learning applied to large-
1319 sample datasets. *Hydrology and Earth System Sciences*, 23, 5089-5110. [https://doi.org/10.5194/hess-23-](https://doi.org/10.5194/hess-23-5089-2019)
1320 [5089-2019](https://doi.org/10.5194/hess-23-5089-2019)
1321

1322 Kratzert, F., Klotz, D., Hochreiter, S., & Nearing, G. S. (2021). A note on leveraging in multiple
1323 meteorological data sets with deep learning for rainfall-runoff modeling. *Hydrology and Earth System*
1324 *Sciences*, 25(5), 2685–2703. <https://doi.org/10.5194/hess-25-2685-2021>.
1325

1326 [Kratzert, F., Nearing, G., Addor, N. et al. \(2023\), Caravan - A global community dataset for large-sample](https://doi.org/10.1038/s41597-023-01975-w)
1327 [hydrology. *Sci Data* 10, 61. https://doi.org/10.1038/s41597-023-01975-w](https://doi.org/10.1038/s41597-023-01975-w)
1328

1329 [Krøgli, I. K., Devoli, G., Colleuille, H., Boje, S., Sund, M., and Engen, I. K.: The Norwegian forecasting](https://doi.org/10.5194/nhess-18-1427-2018)
1330 [and warning service for rainfall- and snowmelt-induced landslides, *Nat. Hazards Earth Syst. Sci.*, 18,](https://doi.org/10.5194/nhess-18-1427-2018)
1331 [1427–1450, https://doi.org/10.5194/nhess-18-1427-2018, 2018.](https://doi.org/10.5194/nhess-18-1427-2018)
1332

1333 [Krysanova, V., Donnelly, C., Gelfan, A., Gerten, D., Arheimer, B., Hattermann, F. and Kundzewicz](https://doi.org/10.1080/02626667.2018.1446214)
1334 [Z.W. \(2018\) How the performance of hydrological models relates to credibility of projections under](https://doi.org/10.1080/02626667.2018.1446214)
1335 [climate change, *Hydrological Sciences Journal*, 63:5, 696-720, DOI: 10.1080/02626667.2018.1446214](https://doi.org/10.1080/02626667.2018.1446214)

1336 Lai, C., Chen, X., Zhong, R., and Wang, Z. (2022), Implication of climate variable selections on the
1337 uncertainty of reference crop evapotranspiration projections propagated from climate variables
1338 projections under climate change, *Agricultural Water Management*, 259(1), 107273,
1339 <https://doi.org/10.1016/j.agwat.2021.107273>.
1340

1341 Lee, D., Lee, G., Kim, S., & Jung, S. (2020). Future Runoff Analysis in the Mekong River Basin under a
1342 Climate Change Scenario Using Deep Learning. *Water*, 12(6):1556. <https://doi.org/10.3390/w12061556>
1343

1344 [Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., et al. \(2021\). Hydrological concept](https://doi.org/10.5194/hess-26-3079-2022)
1345 [formation inside long short term memory \(LSTM\) networks. *Hydrology and Earth System Sciences*, 26](https://doi.org/10.5194/hess-26-3079-2022)
1346 [\(12\), https://doi.org/10.5194/hess-26-3079-2022.](https://doi.org/10.5194/hess-26-3079-2022)
1347

1348 [Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve, P., Slater,](https://doi.org/10.5194/hess-26-3079-2022)
1349 [L., and Dadson, S. J. \(2022\). Hydrological concept formation inside long short-term memory \(LSTM\)](https://doi.org/10.5194/hess-26-3079-2022)
1350 [networks, *Hydrol. Earth Syst. Sci.*, 26, 3079–3101, https://doi.org/10.5194/hess-26-3079-2022.](https://doi.org/10.5194/hess-26-3079-2022)
1351

1352 [Lehner, F., Wahl, E., R., Wood, A. W., Blatchford, D. B., & Llewellyn, D. \(2017\). Assessing recent](https://doi.org/10.1002/2017GL073253)
1353 [declines in Upper Rio Grande runoff efficiency from a paleoclimate perspective. *Geophysical Research*](https://doi.org/10.1002/2017GL073253)
1354 [*Letters*, 44, 4124–4133. https://doi.org/10.1002/2017GL073253](https://doi.org/10.1002/2017GL073253)
1355

1356 Lehner, B., Verdin, K., and Jarvis, A. (2008). New Global Hydrography Derived From Spaceborne
1357 Elevation Data, *Eos T. Am. Geophys. Un.*, 89, 93–94.
1358

1359 Lemaitre-Basset, T., Oudin, L., Thirel, G., and Collet, L.: Unraveling the contribution of potential
1360 evaporation formulation to uncertainty under climate change, *Hydrol. Earth Syst. Sci.*, 26, 2147–2159,
1361 <https://doi.org/10.5194/hess-26-2147-2022>, 2022.
1362

1363 Li, K., Huang, G., Wang, S., Razavi, S., & Zhang, X. (2022). Development of a joint probabilistic
1364 rainfall-runoff model for high-to-extreme flow projections under changing climatic conditions. *Water
1365 Resources Research*, 58, e2021WR031557. <https://doi.org/10.1029/2021WR031557>
1366

1367 Lin, L., Gettelman, A., Fu, Q. et al. Simulated differences in 21st century aridity due to different scenarios
1368 of greenhouse gases and aerosols. *Climatic Change* 146, 407–422 (2018). [https://doi.org/10.1007/s10584-
1369 016-1615-3](https://doi.org/10.1007/s10584-016-1615-3)
1370

1371 [Lin, C., Jain, S., Kim, H., Bar-Joseph, Z. \(2017\). Using neural networks for reducing the dimensions of
1372 single-cell RNA-Seq data, *Nucleic Acids Research*, Volume 45, Issue 17, 29 September 2017, Page e156,
1373 <https://doi.org/10.1093/nar/gkx681>
1374](https://doi.org/10.1093/nar/gkx681)

1375 [Liu, J., Hu, Z., Cui, P., Li, B., and Shen, Z. \(2021\). Heterogeneous risk minimization. In *ICML, PMLR.
1376 PMLR.*
1377](https://proceedings.mlr.press/v97/liu20a.html)

1378

1379 Liu, X., Li, C., Zhao, T., and Han, L. (2020) Future changes of global potential evapotranspiration
1380 simulated from CMIP5 to CMIP6 models, *Atmospheric and Oceanic Science Letters*, 13:6, 568-
1381 575, DOI: 10.1080/16742834.2020.1824983
1382

1383 [Liu, Z., Han, J., and Yang, H. \(2022\), Assessing the ability of potential evaporation models to capture the
1384 sensitivity to temperature, *Agricultural and Forest Meteorology*, 317, 108886.
1385 ~~Liu, Z., Han, J., and Yang, H. \(2022\), Assessing the ability of potential evaporation models to capture the
1386 sensitivity to temperature, *Agricultural and Forest Meteorology*, 317, 108886.~~
1387](https://doi.org/10.1016/j.agrfor.2022.108886)

1388

1389 ~~Liu, Z., Wang T., Han, J., Yang, W., & Yang, H. (2022). Decreases in mean annual streamflow and
1390 interannual streamflow variability across snow affected catchments under a warming climate.
1391 *Geophysical Research Letters*, 49(3), e2021GL097442. <https://doi.org/10.1029/2021GL097442>
1392~~

1393 Lofgren, B.M., Hunter, T.S., Wilbarger, J. (2011), Effects of using air temperature as a proxy for potential
1394 evapotranspiration in climate change scenarios of Great Lakes basin hydrology, *Journal of Great Lakes
1395 Research*, 37 (4), 744-752.
1396

1397 Lofgren, B. M., and Rouhana, J. (2016) Physically Plausible Methods for Projecting Changes in Great
1398 Lakes Water Levels under Climate Change Scenarios. *J. Hydrometeorol.*, 17, 2209–
1399 2223, <https://doi.org/10.1175/JHM-D-15-0220.1>.
1400

1401 Lu, D., Konapala, G., Painter, S. L., Kao, S. C., & Gangrade, S. (2021). Streamflow simulation in data-
1402 scarce basins using Bayesian and physics-informed machine learning models. *Journal of
1403 Hydrometeorology*, 22(6), 1421– 1438. <https://doi.org/10.1175/JHM-D-20-0082.1>
1404

1405

1406 [Lu, J., Sun, G., McNulty, S.G. and Amatya, D.M. \(2005\). A comparison of six potential
1407 evapotranspiration methods for regional use in the southeastern United States. *JAWRA Journal of the
1408 American Water Resources Association*, 41: 621-633. \[https://doi.org/10.1111/j.1752-
1409 1688.2005.tb03759.x\]\(https://doi.org/10.1111/j.1752-

1409 1688.2005.tb03759.x\)](https://doi.org/10.1111/j.1752-1688.2005.tb03759.x)

1410 [Lu, J., Sun, G., McNulty, S.G. and Amatya, D.M. \(2005\). A comparison of six potential](#)
1411 [evapotranspiration methods for regional use in the southeastern United States. JAWRA Journal of the](#)
1412 [American Water Resources Association, 41: 621–633. \[https://doi.org/10.1111/j.1752-\]\(https://doi.org/10.1111/j.1752-1688.2005.tb03759.x\)](#)
1413 [1688.2005.tb03759.x](#)

1414

1415 [Luo, Y., Peng, J. & Ma, J. \(2020\). When causal inference meets deep learning. Nat Mach Intell 2, 426–](#)
1416 [427. <https://doi.org/10.1038/s42256-020-0218-x>](#)

1417

1418 [Ma, J., Yu, M., Fong, S. et al. \(2018\). Using deep learning to model the hierarchical structure and](#)
1419 [function of a cell. Nat Methods 15, 290–298. <https://doi.org/10.1038/nmeth.4627>](#)

1420

1421 Ma, K., Feng, D., Lawson, K., Tsai, W.-P., Liang, C., Huang, X., et al. (2021). Transferring hydrologic
1422 data across continents – leveraging data-rich regions to improve hydrologic prediction in data-sparse
1423 regions. *Water Resources Research*, 57, e2020WR028600. <https://doi.org/10.1029/2020WR028600>

1424

1425 Mai et al. (2022). The Great Lakes runoff intercomparison project phase 4: the Great Lakes (GRIP-GL),
1426 *Hydrologic and Earth System Sciences*, 26 (13), 3537–3573, <https://doi.org/10.5194/hess-26-3537-2022>.
1427

1428 Martens, B., Miralles, D. G., Lievens, H., van der Schalie, R., de Jeu, R. A. M., Fernández-Prieto, D.,
1429 Beck, H. E., Dorigo, W. A., and Verhoest, N. E. C. (2017). GLEAM v3: satellite-based land evaporation
1430 and root-zone soil moisture, *Geosci. Model Dev.*, 10, 1903– 1925, [https://doi.org/10.5194/gmd-10-1903-](https://doi.org/10.5194/gmd-10-1903-2017)
1431 [2017](#).

1432

1433 [Martin, J. T., Pederson, G. T., Woodhouse, C. A., Cook, E. R., McCabe, G. J., Anchukaitis, K. J., et al.](#)
1434 [\(2020\). Increased drought severity tracks warming in the United States’ largest river basin. *Proceedings*](#)
1435 [of the National Academy of Sciences, 117\(21\). <https://doi.org/10.1073/pnas.1916208117>](#)

1436

1437 [Maraun, D. \(2016\). Bias Correcting Climate Change Simulations - a Critical Review. *Curr Clim Change*](#)
1438 [Rep 2, 211–220. <https://doi.org/10.1007/s40641-016-0050-x>](#)

1439

1440 [McCabe, G. J., Wolock, D. M., Pederson, G. T., Woodhouse, C. A., & McAfee, S. \(2017\). Evidence that](#)
1441 [recent warming is reducing upper Colorado River flows. *Earth Interactions*, 21\(10\), 1–14.](#)
1442 [<https://doi.org/10.1175/EI-D-17-0007.1>](#)

1443

1444 [Melsen, L. A., Addor, N., Mizukami, N., Newman, A. J., Torfs, P. J. J. F., Clark, M. P., Uijlenhoet, R., and](#)
1445 [Teuling, A. J. \(2018\). Mapping \(dis\)agreement in hydrologic projections, *Hydrol. Earth Syst. Sci.*, 22,](#)
1446 [1775–1791, <https://doi.org/10.5194/hess-22-1775-2018>.](#)

1447

1448 [Merz, R., Parajka, J., and Blöschl, G. \(2011\), Time stability of catchment model parameters: Implications](#)
1449 [for climate impact analyses, *Water Resour. Res.*, 47, W02531, doi:10.1029/2010WR009505.](#)

1450

1451 Milly, P.C.D. and Dunne, Krista A. (2017). A Hydrologic Drying Bias in Water-Resource Impact
1452 Analyses of Anthropogenic Climate Change. *Journal of the American Water Resources*
1453 *Association (JAWRA)* 53(4): 822– 838. <https://doi.org/10.1111/1752-1688.12538>

1454

1455 [Milly, P. C. D., & Dunne, K. A. \(2020\). Colorado River flow dwindles as warming driven loss of](#)
1456 [reflective snow energizes evaporation. *Science*, 367\(6483\), 1252–1255.](#)
1457 [<https://doi.org/10.1126/science.aay9187>](#)

1458

1459 [Monteith, J. L. \(1965\), Evaporation and environment, in: *Symposia of the society for experimental*](#)
1460 [biology, volume 19, Cambridge University Press \(CUP\), Cambridge, UK, 205–234 pp.](#)

1461
1462 Mote, P. W., Li, S., Lettenmaier, D. P., Xiao, M., & Engel, R. (2018). Dramatic declines in snowpack in
1463 the western US. *npj Climate and Atmospheric Science*, 1:2. <https://doi.org/10.1038/s41612-018-0012-1>
1464
1465 NACLMS: NACLMS website, [http://www.cec.org/north-american-environmental-atlas/land-cover-2010-](http://www.cec.org/north-american-environmental-atlas/land-cover-2010-landsat-30m/)
1466 [landsat-30m/](http://www.cec.org/north-american-environmental-atlas/land-cover-2010-landsat-30m/) (last access: 31 May 2023), 2017.
1467
1468 [Najibi, N., Mukhopadhyay, S., & Steinschneider, S. \(2022\). Precipitation scaling with temperature in the](#)
1469 [Northeast US: Variations by weather regime, season, and precipitation intensity. *Geophysical Research*](#)
1470 [Letters, 49, e2021GL097100. <https://doi.org/10.1029/2021GL097100>](#)
1471
1472 Nash, J. E. and Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I – A
1473 discussion of principles, *J. Hydrol.*, 10, 282–290.
1474
1475 Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., et al. (2021). What
1476 role does hydrological science play in the age of machine learning? *Water Resources Research*, 57,
1477 e2020WR028091. <https://doi.org/10.1029/2020WR028091>
1478
1479 [Nearing, G. S., Klotz, D., Frame, J. M., Gauch, M., Gilon, O., Kratzert, F., Sampson, A. K., Shalev, G.,](#)
1480 [and Nevo, S. \(2022\). Technical note: Data assimilation and autoregression for using near-real-time](#)
1481 [streamflow observations in long short-term memory networks. *Hydrol. Earth Syst. Sci.*, 26, 5493–5513,](#)
1482 [https://doi.org/10.5194/hess-26-5493-2022.](#)
1483
1484 Newman, A., Clark, M. P., Sampson, K., Wood, A., Hay, L., Bock, A., et al. (2015). Development of a
1485 large-sample watershed-scale hydrometeorological dataset for the contiguous USA: Data set
1486 characteristics and assessment of regional variability in hydrologic model performance. *Hydrology and*
1487 *Earth System Sciences*, 19(1), 209-223. <https://doi.org/10.5194/hess-19-209-2015>
1488
1489 Nordling, K., Korhonen, H., Raisanen, J., Partanen, A.-I., Samset, B.H., and Merikanto, J. (2021),
1490 Understanding the surface temperature response and its uncertainty to CO₂, CH₄, black carbon, and
1491 sulfate, *Atmos. Chem. Phys.*, 21, 14941-14958.
1492
1493 [Olsson, J., and Lindstrom, G. \(2008\), Evaluation and calibration of operational hydrological ensemble](#)
1494 [forecasts in Sweden *Journal of Hydrology*, 350 \(1–2\), 14-24.](#)
1495
1496 [Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andreassian, V., Anctil, F., and Loumagne,](#)
1497 [C. \(2005\). Which potential evapotranspiration input for a lumped rainfall-runoff model? Part 2—Towards](#)
1498 [a simple and efficient potential evapotranspiration model for rainfall-runoff modeling. *Journal of*](#)
1499 [Hydrology 303: 290–306.](#)
1500
1501 [Plesca, I., Timbe, E., Exbrayat, J.F., Windhorst, D., Kraft, P., Crespo, P., Vachéa, K.B., Frede, H.G., and](#)
1502 [Breuer, L. \(2012\). Model intercomparison to explore catchment functioning: Results from a remote](#)
1503 [montane tropical rainforest. *Ecol. Model.*, 239, 3–13.](#)
1504
1505 Priestley, C. H. B., and Taylor, R. J. (1972). On the Assessment of Surface Heat Flux and Evaporation
1506 Using Large-Scale Parameters. *Mon. Wea. Rev.*, 100, 81–92, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0493(1972)100<0081:OTAOSH>2.3.CO;2)
1507 [0493\(1972\)100<0081:OTAOSH>2.3.CO;2](https://doi.org/10.1175/1520-0493(1972)100<0081:OTAOSH>2.3.CO;2).
1508
1509 Pryor, S.C., Barthelmie, R.J., Bukovsky, M.S. et al. Climate change impacts on wind power
1510 generation. *Nat Rev Earth Environ* 1, 627–643 (2020). <https://doi.org/10.1038/s43017-020-0101-7>

1511
1512 Razavi, S. (2021). Deep learning, explained: Fundamentals, explainability, and bridgeability to process-
1513 based modelling, *Environmental Modelling and Software*,
1514 105159, <https://doi.org/10.1016/j.envsoft.2021.105159>.
1515
1516 [Reichert, P., Ma, K., Höge, M., Fenicia, F., Baity-Jesi, M., Feng, D., and Shen, C.: Metamorphic Testing](#)
1517 [of Machine Learning and Conceptual Hydrologic Models, *Hydrol. Earth Syst. Sci. Discuss.* \[preprint\],](#)
1518 <https://doi.org/10.5194/hess-2023-168>, in review, 2023.
1519
1520 [Rungee, J., Ma, Q., Goulden, M. L., & Bales, R. \(2021\). Evapotranspiration and runoff patterns across](#)
1521 [California’s Sierra Nevada. *Frontiers in Water*, 3:655485. <https://doi.org/10.3389/frwa.2021.655485>](#)
1522
1523 [Safeeq, M., Bart, R. R., Pelak, N. F., Singh, C. K., Dralle, D. N., Hartsough, P., & Wagenbrenner, J. W.](#)
1524 [\(2021\). How realistic are water-balance closure assumptions? A demonstration from the southern sierra](#)
1525 [critical zone observatory and kings river experimental watersheds. *Hydrological Processes*, 35: e14199.](#)
1526 <https://doi.org/10.1002/hyp.14199>
1527
1528 [Seibert, J. and Bergström, S. \(2022\). A retrospective on hydrological catchment modelling based on half a](#)
1529 [century with the HBV model, *Hydrol. Earth Syst. Sci.*, 26, 1371–1388, \[https://doi.org/10.5194/hess-26-\]\(https://doi.org/10.5194/hess-26-1371-2022\)](#)
1530 [1371-2022](https://doi.org/10.5194/hess-26-1371-2022).
1531
1532 Shangguan, W., Dai, Y., Duan, Q., Liu, B., and Yuan, H. (2014). A global soil data set for earth system
1533 modeling, *J. Adv. Model. Earth Sy.*, 6, 249–263.
1534
1535 Shaw, S.B. and Riha, S.J. (2011), Assessing temperature-based PET equations under a changing climate
1536 in temperate, deciduous forests. *Hydrol. Process.*, 25: 1466-1478. <https://doi.org/10.1002/hyp.7913>
1537
1538 [Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., and Cui, P. \(2021\). Towards out-of-distribution](#)
1539 [generalization: A survey. arXiv preprint arXiv:2108.13624.](#)
1540
1541 [Siddik, M.A.B., Dickson, K.E., Rising, J. et al. Interbasin water transfers in the United States and](#)
1542 [Canada. *Sci Data* 10, 27 \(2023\). <https://doi.org/10.1038/s41597-023-01935-4>](#)
1543
1544 Steinman, A.D. et al. (2017), Ecosystem services in the Great Lakes, *Journal of Great Lakes Research*, 43
1545 (3), 161-168. <https://doi.org/10.1016/j.jglr.2017.02.004>
1546
1547 [Stewart, I. T., Cayan, D. R., & Dettinger, M. D. \(2005\). Changes toward Earlier Streamflow Timing](#)
1548 [across Western North America, *Journal of Climate*, 18\(8\), 1136–1155.](#)
1549 <https://doi.org/10.1175/JCLI3321.1>
1550
1551 [Su, Q., & Singh, V. P. \(2023\). Calibration-free Priestley-Taylor method for reference evapotranspiration](#)
1552 [estimation. *Water Resources Research*, 59, e2022WR033198. <https://doi.org/10.1029/2022WR033198>](#)
1553
1554 Szilagyi, J., Crago, R., and Qualls, R. (2017), A calibration-free formulation of the complementary
1555 relationship of evaporation for continental-scale hydrology, *J. Geophys. Res. Atmos.*, 122, 264– 278,
1556 doi:10.1002/2016JD025611.
1557
1558

1559 [Taranu, I.S., Somot, S., Alias, A. et al. Mechanisms behind large-scale inconsistencies between](#)
1560 [regional and global climate model-based projections over Europe. *Clim Dyn* 60, 3813–3838](#)
1561 [\(2023\). <https://doi.org/10.1007/s00382-022-06540-6>](#)

1562
1563 [Towler, E., Foks, S. S., Dugger, A. L., Dickinson, J. E., Essaid, H. I., Gochis, D., Viger, R. J., and Zhang,](#)
1564 [Y. \(2023\): Benchmarking high-resolution hydrologic model performance of long-term retrospective](#)
1565 [streamflow simulations in the contiguous United States, *Hydrol. Earth Syst. Sci.*, 27, 1809–1825,](#)
1566 <https://doi.org/10.5194/hess-27-1809-2023>.

1567
1568 [Vasudevan, R.K., Ziatdinov, M., Vlcek, L. et al. \(2021\). Off-the-shelf deep learning is not enough, and](#)
1569 [requires parsimony, Bayesianity, and causality. *npj Comput Mater* 7, 16. \[https://doi.org/10.1038/s41524-\]\(https://doi.org/10.1038/s41524-020-00487-0\)](#)
1570 [020-00487-0](https://doi.org/10.1038/s41524-020-00487-0)

1571
1572 [Wallner, M., and Haberlandt, U. \(2015\). Non-stationary hydrological model parameters: a framework](#)
1573 [based on SOM-B. *Hydrol. Process.*, 29, 3145–3161. doi: 10.1002/hyp.10430.](#)

1574
1575 Wang, Q. J. (1991). The genetic algorithm and its application to calibrating conceptual rainfall-runoff
1576 models, *Water Resources Research*, 27(9), 2467-2471. <https://doi.org/10.1029/91WR01305>

1577
1578 [Wang, J., Lan, C., Liu, C., Ouyang, Y., Qin, T., Lu, W., Chen, Y., Zeng, W., Yu, P.S.](#)
1579 [\(2023\). Generalizing to Unseen Domains: A Survey on Domain Generalization, in *IEEE Transactions on*](#)
1580 [Knowledge and Data Engineering](#), vol. 35, no. 8, pp. 8052-8072, 1 Aug. 2023, doi:
1581 [10.1109/TKDE.2022.3178128](https://doi.org/10.1109/TKDE.2022.3178128).

1582
1583 Wi, S., & Steinschneider, S. (2022). Assessing the physical realism of deep learning hydrologic model
1584 projections under climate change. *Water Resources Research*, 58,
1585 e2022WR032123. <https://doi.org/10.1029/2022WR032123>

1586
1587 [Wolock, D.M., McCabe, G.J. \(1999\). Estimates of runoff using water balance and atmospheric general](#)
1588 [circulation models. *Journal of the American Water Resources Association* 35: 1341–1350.](#)

1589
1590 [Woodhouse, C. A., & Pederson, G. T. \(2018\). Investigating runoff efficiency in upper Colorado river](#)
1591 [streamflow over past centuries. *Water Resources Research*, 54, 286–300.](#)
1592 <https://doi.org/10.1002/2017WR021663>

1593
1594 [Wu, H., Zhu, W., and Huang, B. \(2021\), Seasonal variation of evapotranspiration, Priestley-Taylor](#)
1595 [coefficient and crop coefficient in diverse landscapes, *Geography and Sustainability*, 2\(3\), 224-233,](#)
1596 <https://doi.org/10.1016/j.geosus.2021.09.002>

1597
1598 [Yan, H., Sun, N., Eldardiry, H., Thurber, T. B., Reed, P. M., Malek, K., et al. \(2023\). Large ensemble](#)
1599 [diagnostic evaluation of hydrologic parameter uncertainty in the Community Land Model Version 5](#)
1600 [\(CLM5\). *Journal of Advances in Modeling Earth Systems*, 15,](#)
1601 [e2022MS003312. <https://doi.org/10.1029/2022MS003312>](https://doi.org/10.1029/2022MS003312)

1602
1603 [Yang, Y., & Chui, T. F. M. \(2021\). Reliability assessment of machine learning models in hydrological](#)
1604 [predictions through metamorphic testing. *Water Resources Research*, 57,](#)
1605 [e2020WR029471. <https://doi.org/10.1029/2020WR029471>](https://doi.org/10.1029/2020WR029471)

1606

1607 Yilmaz, K. K., Gupta, H. V., and Wagener, T. (2008). A process-based diagnostic approach to model
1608 evaluation: Application to the NWS distributed hydrologic model, *Water Resour. Res.*, 44, 1–18.

1609
1610 Zhong, L., Lei, H., & Gao, B. (2023). Developing a physics-informed deep learning model to simulate
1611 runoff response to climate change in Alpine catchments. *Water Resources Research*, 59,
1612 e2022WR034118. <https://doi.org/10.1029/2022WR034118>

1613
1614

1615

1616

1617

1618

1619

1620

1621

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

Supplemental Material for

**On the need for physical constraints in deep learning rainfall-runoff
projections under climate change: [a sensitivity analysis to warming and shifts
in potential evapotranspiration](#)**

Sungwook Wi¹, Scott Steinschneider¹

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

Summary

This supplementary material file contains [one supplemental section of text, ~~six eight figures, one supplemental section of text, and three additional tables, and eightnine figures~~](#) in support of the analysis and conclusions presented in the main article.

Text S1: Adjustments to Static Attributes

In the primary article, we describe two sets of scenarios for the deep learning models used in this work: 1) one in which changes are only made to the dynamic inputs features of each model, and 2) one with changes to both dynamic features and to static features that depend on those dynamic features. Here we describe in more detail the adjustments made to the static features for each site, which include: pet_mean, aridity, t_mean, frac_snow (see Table S1 below for the definition of these features). Importantly, these are the static features that are dependent on temperature and PET, the two dynamic inputs adjusted in our analysis.

To adjust t_mean, we use the full time series of daily maximum and minimum temperature (on which t_mean was originally based), and shift those time series upward by 4°C. Using those adjusted series, we calculate daily average temperature as the mean of maximum and minimum temperature on each day, and then calculate the long-term mean of daily average temperature to develop an updated estimate of t_mean.

To adjust frac_snow, we first calculate the adjusted time series of daily average temperature based on the time series of daily maximum and minimum temperature shifted upward by 4°C. Then, we count all days in the record when precipitation occurs and this adjusted time series of daily average temperature is below 0°C, and divide this number by the total number of days of non-zero precipitation in the record. The resulting value is the updated value for frac_snow.

We develop two versions of adjusted pet_mean, one based on Hamon PET and the other for Priestley-Taylor PET. The adjusted Hamon PET is based entirely on the series of daily maximum and minimum temperature shifted by 4°C. We use Eqs. 7-8 in the main article to calculate daily Hamon PET under warming. We then take the long-term mean of this time series to develop an updated estimate of pet_mean. Similarly, for Priestley-Taylor PET, we couple the warmed temperature time series with the unadjusted time series of net shortwave radiation, and then use the approach in Eq. 9 in the main article to calculate a daily time series of Priestley-Taylor PET. We again take the long-term mean of this time series to develop an updated estimate of pet_mean.

Finally, we develop two versions of adjusted aridity, one based on Hamon PET and the other for Priestley-Taylor PET. In both cases, we calculate adjusted aridity as the ratio of the updated values for pet_mean under warming and the unadjusted value for long-term mean precipitation (another static input to the models).

Table S1. Static watershed attributes that are adjusted in a subset of scenarios used in this analysis.

<u>Attribute</u>	<u>Description</u>
<u>pet_mean</u>	<u>Mean daily potential evapotranspiration</u>
<u>aridity</u>	<u>Ratio of mean PET to mean precipitation</u>
<u>t_mean</u>	<u>Mean of daily maximum and daily minimum temperature</u>

<u>frac_snow</u>	<u>Fraction of precipitation falling on days with mean daily temperatures below 0°C</u>
----------------------------------	---

[Additional Supporting Tables](#)

[Table S2. Range of values considered in the grid search during hyper-parameter tuning.](#)

<u>Hyper-parameter</u>	<u>Values Tested</u>
<u>Number of Hidden Layer Nodes</u>	<u>64, 96, 128, 256</u>
<u>Mini-Batch Size</u>	<u>64, 128, 256, 512</u>
<u>Learning Rate</u>	<u>0.0001, 0.0005, 0.001, 0.005</u>
<u>Number of Epochs</u>	<u>30, 50</u>
<u>Dropout Rate*</u>	<u>0, 0.2, 0.4</u>

[Table S3. Additional details for gauges highlighted in Figures 5 and 6 of main article.](#)

<u>Gauge ID</u>	<u>Country</u>	<u>Site Name</u>	<u>Drainage Area (km²)</u>
<u>02ED032</u>	<u>Canada</u>	<u>Willow Creek near Minesing</u>	<u>231</u>
<u>02GG013</u>	<u>Canada</u>	<u>Black Creek near Bradshaw</u>	<u>213</u>
<u>02HJ003</u>	<u>Canada</u>	<u>Ouse River near Westwook</u>	<u>283</u>
<u>04126740</u>	<u>United States</u>	<u>Platte River at Honor, MI</u>	<u>324</u>
<u>04220045</u>	<u>United States</u>	<u>Oak Orchard Creek near Shelby NY</u>	<u>378</u>
<u>04168400</u>	<u>United States</u>	<u>Lower River Rouge at Dearborn, MI</u>	<u>236</u>

[Additional Supporting Figures](#)

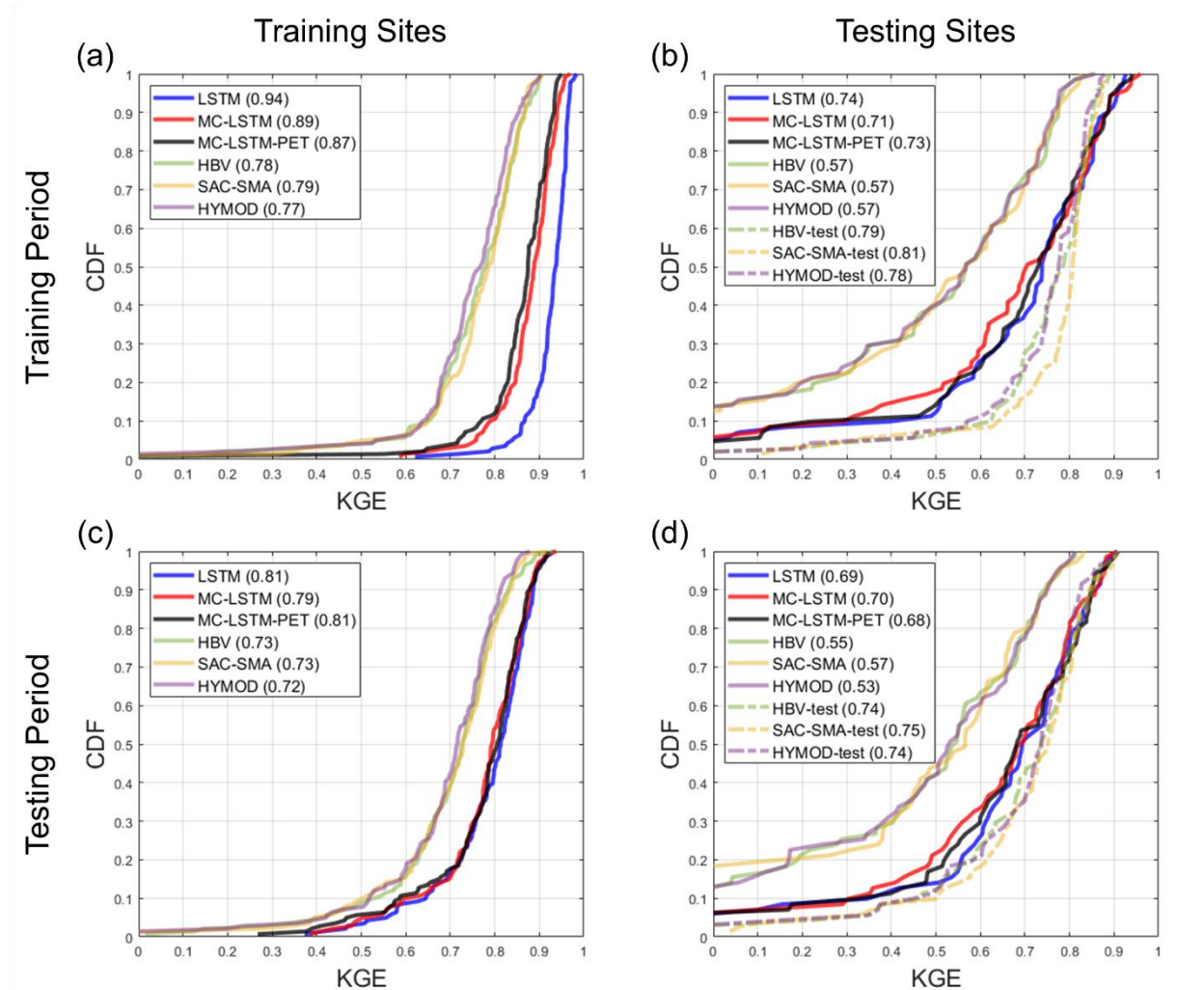


Figure S1. The distribution of Kling-Gupta efficiency (KGE) for streamflow estimates across sites from each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results for the testing period are shown in panels (c) and (d), respectively. For the process models fit to the testing sites (denoted “-test”), no performance results are available at the training sites. All models are trained using Hamon PET.

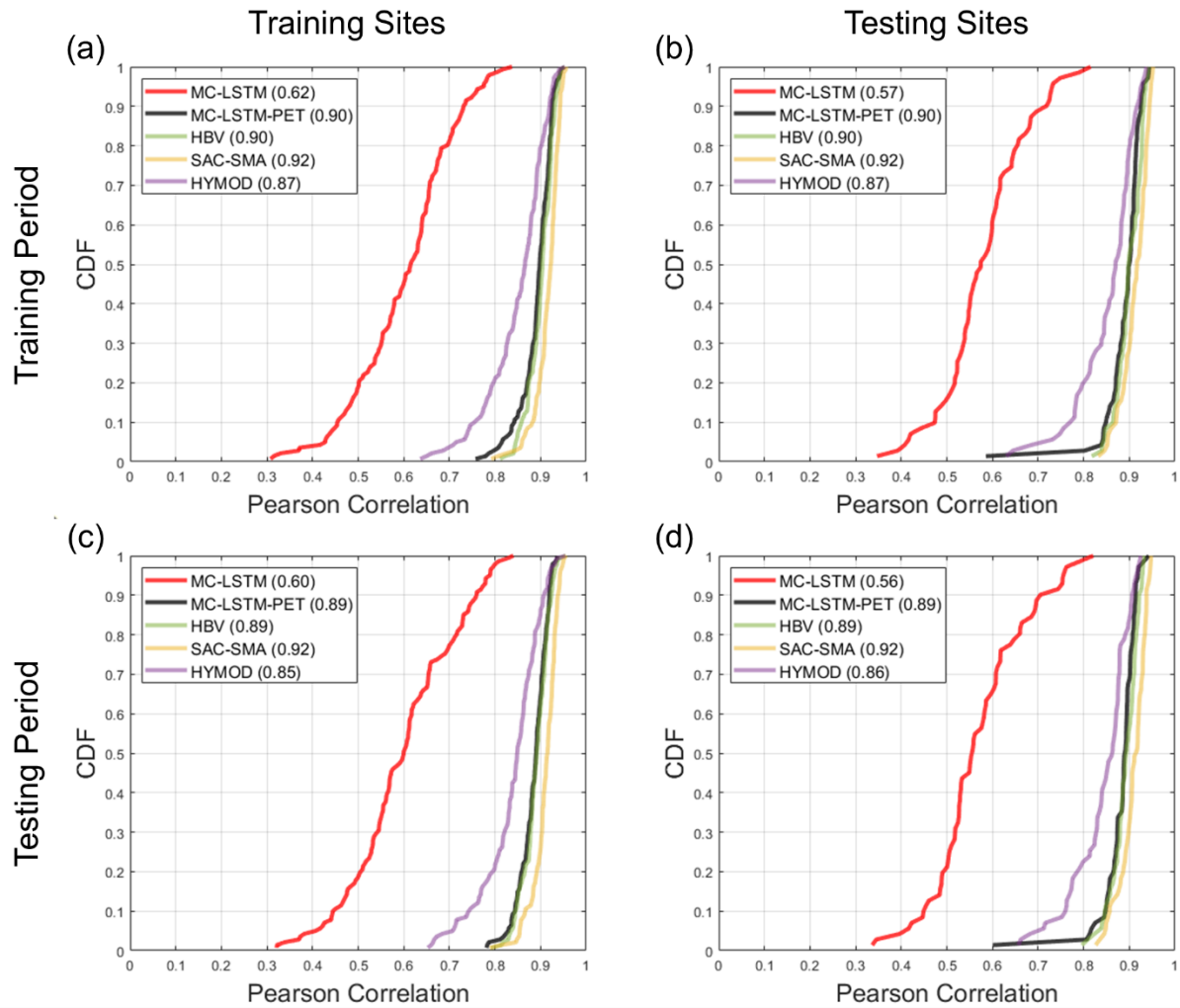


Figure S2. The correlation between model estimated and [observed-GLEAM](#) AET from each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results for the testing period are shown in panels (c) and (d), respectively. The LSTM is not included in this comparison. All models are trained using Priestley-Taylor PET.

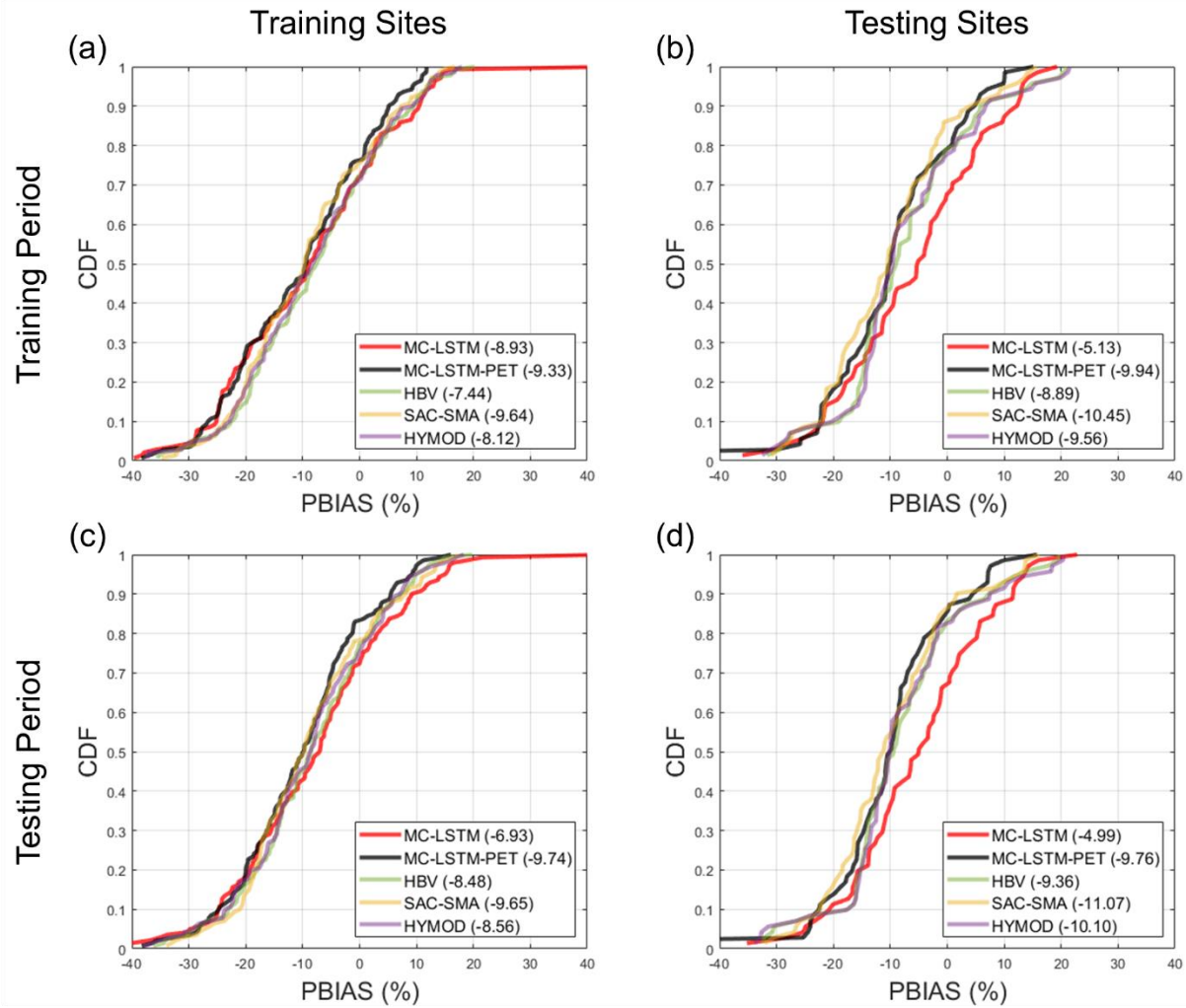


Figure S3. The PBIAS between model estimated and [GLEAM observed](#) AET from each model at the (a) the 141 training sites and (b) 71 testing sites for the training period. Similar results for the testing period are shown in panels (c) and (d), respectively. The LSTM is not included in this comparison. All models are trained using Priestley-Taylor PET.

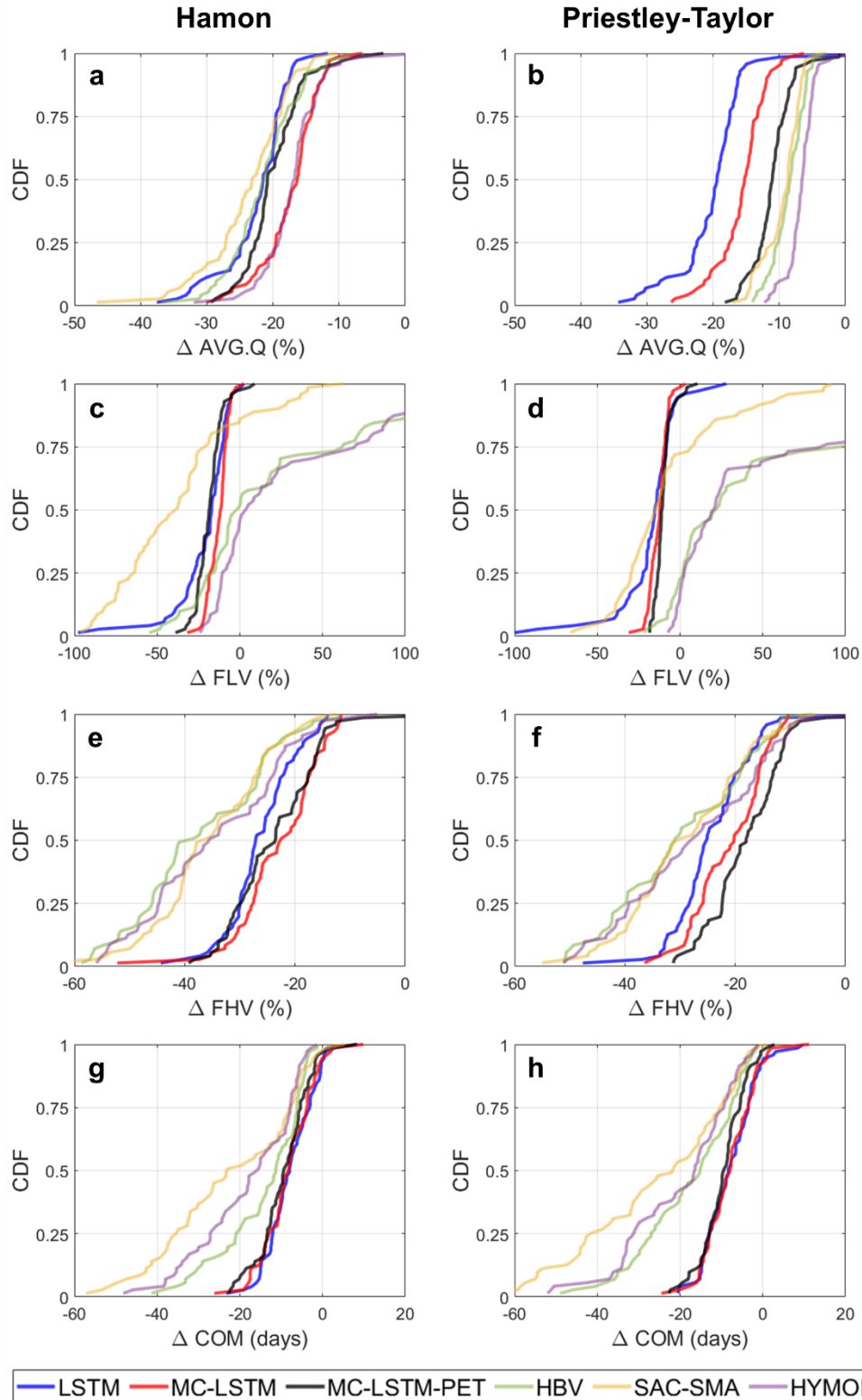


Figure S45. The distribution of change in (a,b) AVG.Q, (c,d) FLV, (e,f) FHV, and (g,h) COM across the 71 testing sites and all models under a scenario of 4°C warming using (a,c,e,g) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the DL models, changes were only made to the dynamic inputs (i.e., no changes to static inputs).

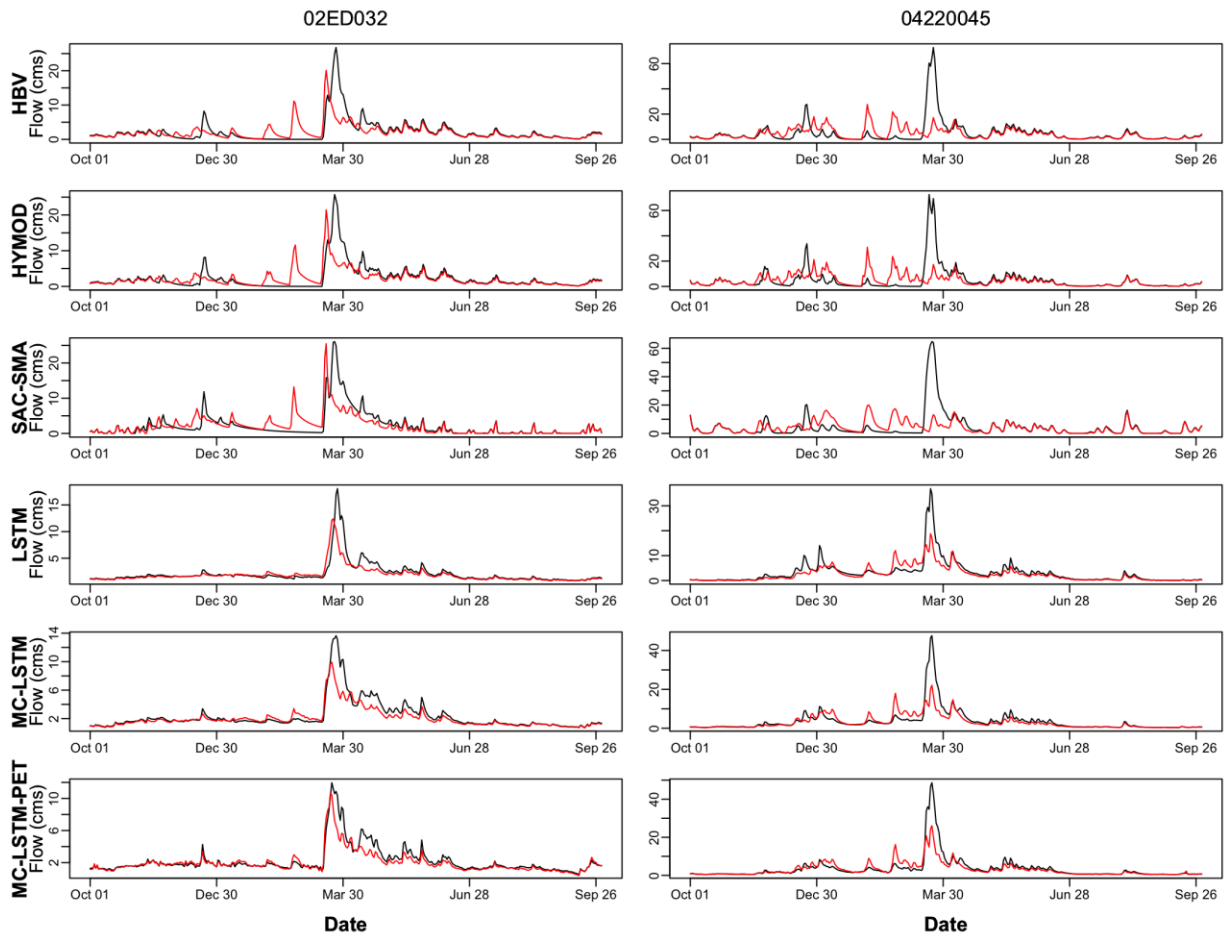


Figure S5. Daily streamflow hydrograph for one water year (2002 October- 2003 September) across the three different process-based models (HBV, HYMOD, SAC-SMA) and deep-learning models (LSTM, MC-LSTM, MC-LSTM-PET) under 0°C warming (black) and 4°C warming (red). Results are shown for two sites (highlighted in Figure 1 of the main article), and are constructed with models using Priestley-Taylor PET.

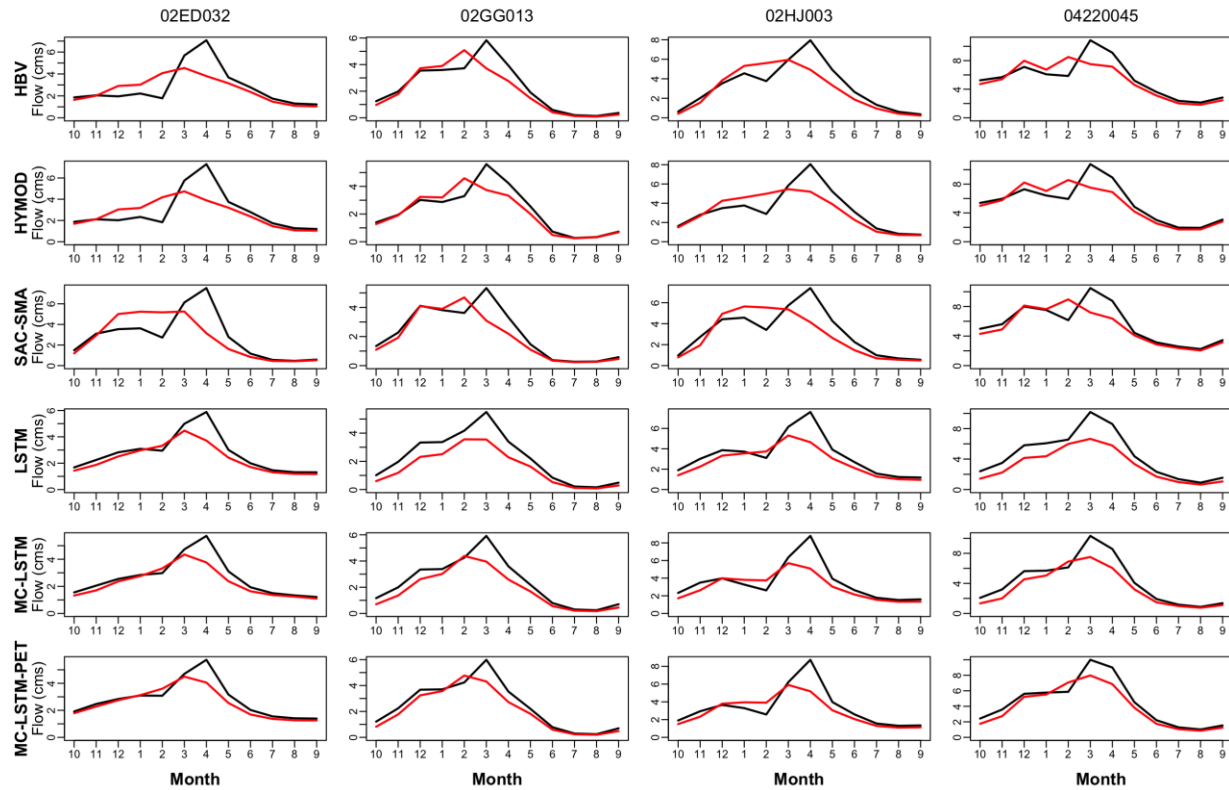


Figure S6. Mean monthly streamflow averaged across the entire record, shown throughout the water year (October-September) across for the three different process-based models (HBV, HYMOD, SAC-SMA) and deep-learning models (LSTM, MC-LSTM, MC-LSTM-PET) under 0°C warming (black) and 4°C warming (red). Results are shown on a water year basis (October-September) for four sites (highlighted in Figure 1 of the main article), and are constructed with models using Priestley-Taylor PET.

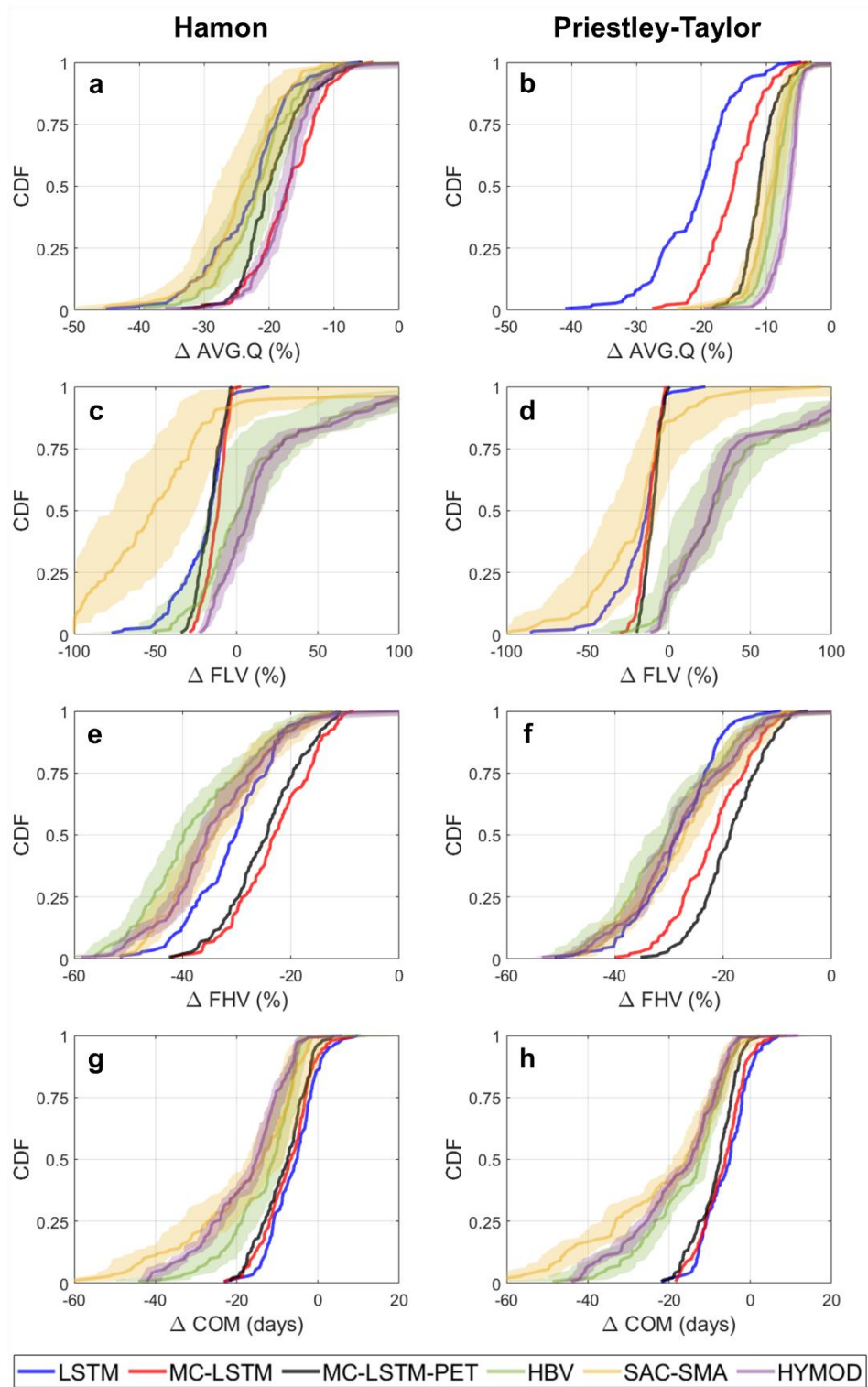


Figure S7. The distribution of change in (a,b) long term mean daily flow (AVG.Q), (c,d) low flows (FLV), (e,f) high flows (FHV), and (g,h) seasonal streamflow timing (COM) across the 141 training sites and all models under a scenario of 4°C warming using (a,c,e,g) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the deep learning models, changes were only made to the dynamic inputs (i.e., no changes to static inputs). For the process models, the uncertainty in the change in each streamflow attribute across 10 different training trails is shown as translucent shading.

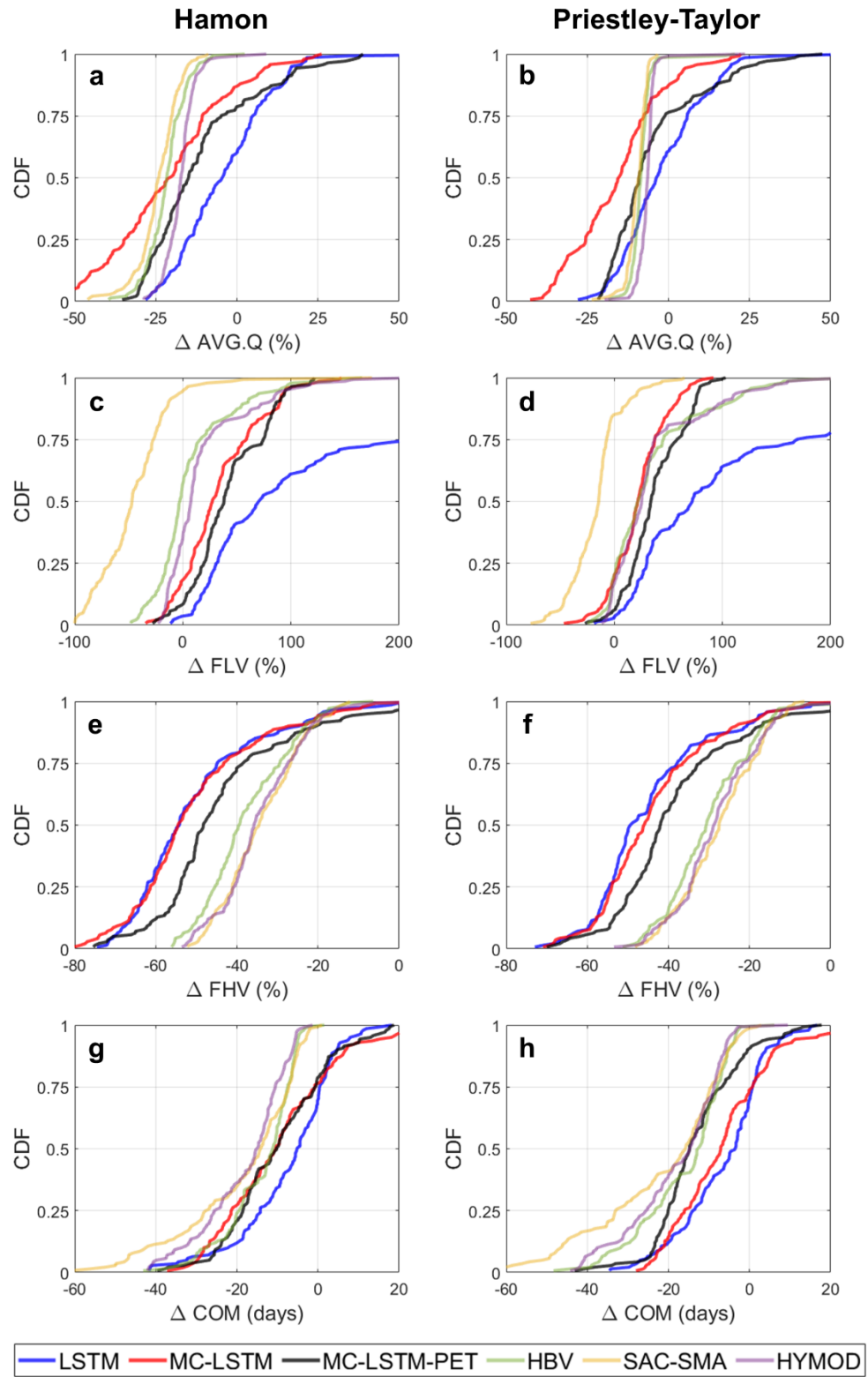


Figure S84. The distribution of change in (a,b) AVG.Q, (c,d) FLV, (e,f) FHV, and (g,h) COM across the 141 training sites and all models under a scenario of 4°C warming using (a,c,e,g) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the DL models, changes were made to both the dynamic and static inputs.

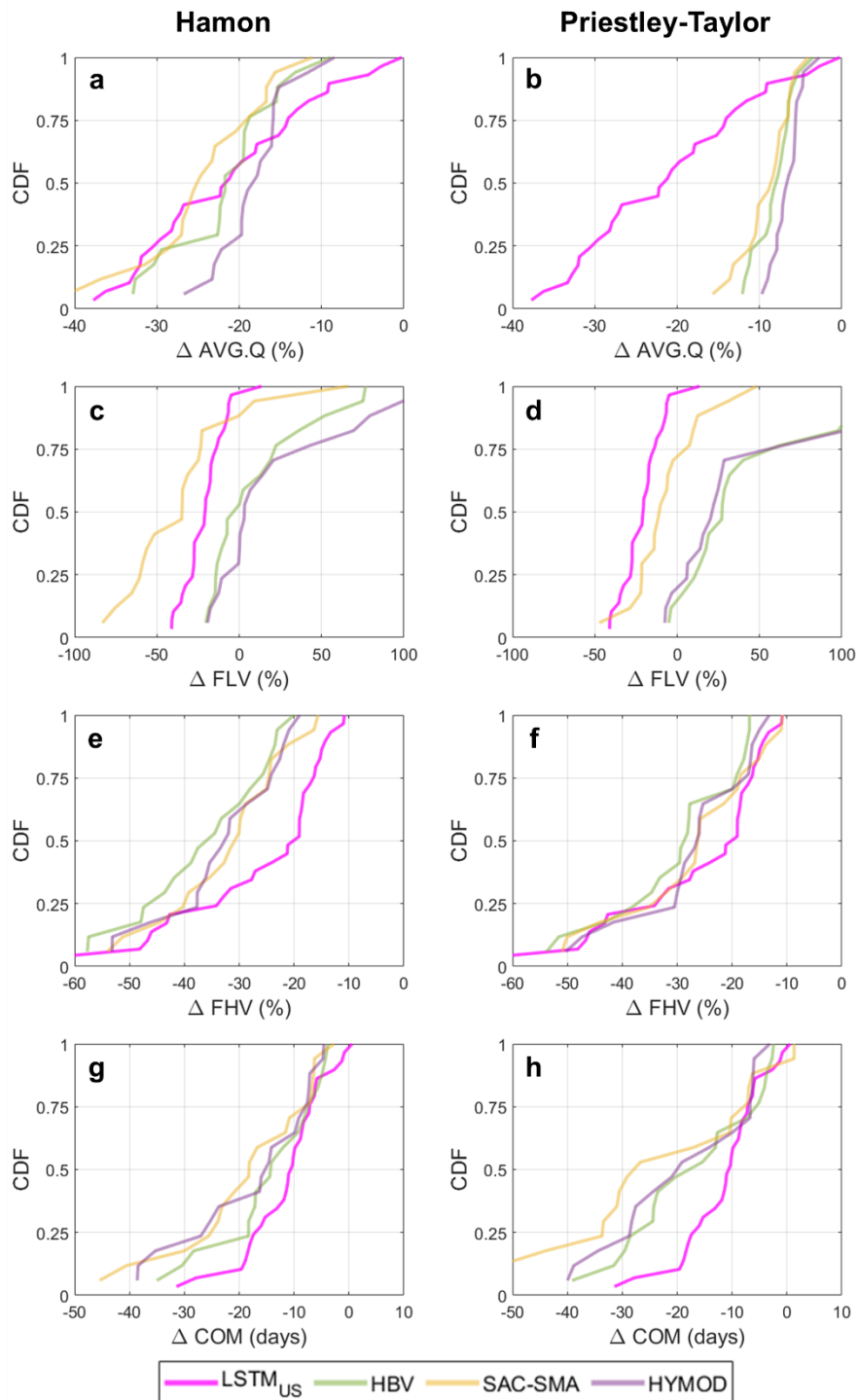


Figure S96. The distribution of change in (a,b) AVG.Q, (c,d) FLV, (e,f) FHV, and (g,h) COM across 29 CAMELS sites within the Great Lakes basin under the National LSTM, as well as for 17 of those 29 sites from the Great Lakes process models, under a scenario of 4°C warming. For the process models only, results differ when using (a,c,e,f) Hamon PET and (b,d,f,h) Priestley-Taylor PET. For the National LSTM, changes were made to both the dynamic and static inputs.