1 2 3 4	Hydrology and Earth System Science Manuscript #HESS-2024-256 January 25 th , 2025	25,
5 6 7 8 9 10 11	Subject: Response to Comments on Rev Temperature Modeling: a review and	view Paper " Machine Learning in River/Stream Water I metrics for evaluation"
12 13 14	Dear Dr. Christa Kelleher, Mr. Jeremy	Diaz and referees #1 and #2,
15 16 17 18 19	We thank you for your time and patience line with referee feedback and think the as our written author response to referee	ce in reviewing our manuscript. We have revised the manuscript in e manuscript is much improved as a result. This document serves e comments.
20 21 22 23	This document separates responses by r $21 - 30$). Where referee comments had both indicate that the part belonged to c all comments were kept in their original	referee: referee #1 (pages 2 - 14), #2 (pages 15 - 20), and #3 (pages several parts, we separated the comments into "1a, 1b, etc.," to one comment, and allow for a more organized response. Otherwise, 1 format.
24 25 26 27 28	For revisions, new/edited text is in BLU statement " <i>revised lines XXX-XXX</i> " ind Track-Changes-Manuscript-HESS-2 manuscript " 3-Clean-Manuscript-HES	JE, while removed text is shown as being crossed out. The icates the in-line placement in the track-changes manuscript "2-024-256", where the described changes will be found. The SS-2024-256" is the "final" version.
29 30 31 32 33	Once again, thank you kindly for your t	ime and consideration.
34 35 36	Sincerely,	
37 38 39 40	Cor Cor	Torrie Schogue
40 41 42 43 44 45	Claudia R. Corona Postdoctoral Fellow Colorado School of Mines	Terri S. Hogue Dean, Earth and Society Programs Colorado School of Mines
46 47 48		
49 50		

1 <u>Referee #1 Comments</u>

2 The manuscript on "ML in Stream/River Water Temperature Modeling, a review and metrics for

evaluation" focuses on providing a comprehensive review of ML studies, including traditional and recent
methods in ML and AI, on stream temperature modeling and prediction. Overall, the manuscript is wellwritten and covers most of the relevant papers, but there are a few strategic points I would like to share

- 6 with the authors:
- 7

8 AUTHOR RESPONSE: We appreciate the referee's feedback and think the manuscript is much

9 improved as a result. For reference, we separated some referee comments into a, b, etc., to provide a more
 10 organized response. Thank you for your time and insight. Proposed new/edited text is in BLUE. Revised

11 lines in the track-changes manuscript are indicated by the statement: (*revised lines XXX-XXX*).

12

13

14 **1a.** Figures 1 & 2 & 3 & table 2: The manuscript provides a table for multiple metrics such as R2, NSE, 15 RMSE, and MAE, and suggested a rate of numbers to rate the ML methods' performances. This table is 16 based on the metrics that have been achieved by the studies in the previous years which are reflected in 17 figures 1 & 2 & 3. However, those studies vary in terms of case studies, number of basins included in the study, running regional or local models. We know that ML models are prone to overfitting, especially for 18 19 stream temperature that follows a relatively sinusoidal curve through a year, which means it is more 20 predictable for complex models such as LSTM. However, it means the models are prone to easily overfit. 21 Therefore, I suggest the authors encourage the stream temperature researchers to go towards making more 22 generalizable models and less overfitted. For example, instead of suggesting performance metrics, the 23 authors can provide a few steps to make sure the models are not overfitted or underfitted. For instance, 24 always considering a spatial test on ungauged sites (basins). We know that spatial tests are more difficult 25 tasks rather than temporal tests.

25 26

AUTHOR RESPONSE: We agree that the SWT studies vary spatially/temporally and that ML models
risk overfitting. We appreciate the referee's comments in pointing out areas of improvement and we
suggest adding the following: 1) a new subsection 2.4.1 "SWT Predictions using ML" on
overfitting/underfitting that suggests the need for temporal- and spatially-focused testing as suggested by
the referee, and 2) a diagram showing initial steps to mitigate overfitting. The new text is below:

32 33

34

47

*new Section 2.4.1, Identifying Model Complexity (revised lines 464-483)

35 The strong success of ML-use in SWT modeling warrants a brief and broad overview on identifying 36 model complexity to minimize overfitting and underfitting" of models. When a model is too complex, 37 i.e., has too many features or parameters relative to the number of observations, or is forced to 38 overextend its capabilities, i.e., make predictions with insufficient training data, the model runs the 39 risk of overfitting (Srivastava et al., 2014). An overfitted model fits the training data "too well", capturing noise and details that provide high accuracy on a training dataset, only to perform poorly 40 once the model encounters "unseen" data in testing/validation (Xu and Liang, 2021). Scenarios where 41 42 overfitting may be temporarily acceptable are: 1) model development is at preliminary stages, the 43 interest is in a "proof of life" concept, 2) when the objective is to identify heavily-relied on features 44 by the model, i.e., feature importance, or 3) in highly-controlled modeling environments where the 45 expected data will be consistently similar to the training dataset. The latter is more likely in industrial 46 applications and unlikely in the changing nature of hydrology.

In contrast, underfitting occurs when a model is too simple to capture any patterns in the data, which
can also lead to unsatisfactory performance in training, testing and validation. Underfitting can occur
with inadequate model features, poor model complexity or when regularization techniques, (e.g., L1
or L2 regularization), are over-used, making the model too rigid and unable to respond to changes in
the data. Given the propensity of ML models to effectively learn the training data, underfitting is less
of an issue in ML whereas overfitting can be widespread. In Figure 1, we present an example

- 1 workflow that researchers can use to transition away from overfitting and towards generalizability. In 2 the five-step outline (Fig. 1), we suggest the need for "Temporal, Unseen, Ungaged Region Tests" 3 (TUURTs), which is a call for temporal and spatially-focused testing that can be used to strengthen 4 model robustness.
 - Revised lines 484-486:

5 6



Figure 1. Diagram outlining steps that can be taken in modeling process to mitigate overfitting.

1b. Therefore, it is acceptable to get lower performance on ungauged basins, however, the metrics should not be vastly different from temporal tests. A more challenging experiment is to test the trained model on 13 regions that have not been seen by the model. In theory, if a model has been able to capture true relations 14 between the driving factors on stream temperature, it should achieve a relatively decent performance on 15 basins with different hydrologic, geologic, and climatic characteristics from the trained basins. As a 16 researcher on SWT, I would rather to have a model that passes all these three tests (temporal, ungaged, 17 unseen regions) with relatively close metrics, rather than having a model that gives high performance in 18 temporal tests and low performance in the other two tests.

19

20 AUTHOR RESPONSE: We agree. The referee mentions a key point that having a SWT model pass all 21 three tests for temporal, ungaged, and unseen regions may be more qualitatively sound, but as of this 22 review, we had not yet seen any ML-SWT papers that test for all three cases. For example, Topp et al. 23 (2023) held out a region to be considered "unseen" but did not test for ungaged basins. Hani et al. (2023) used an inverse weighted distance interpolation method to estimate values for ungaged sites but did not 24 25 test for "unseen" data. Souaissi et al. (2023) used a leave-one-out cross-validation technique to mimic

- 1 arguably not testing for new, ungaged sites but rather "unseen" (i.e., tested only within the development
- 2 dataset, not for new sites). Siegel et al. (2023), a non-ML paper tested for "ungaged" and "unseen" data,
- 3 but did not perform a temporal test. A newly published example, Philippus et al. (2024), which
- 4 considered spatial testing on ungaged basins, has been added. We further agree with the theory posited
- by the referee that a model capturing true relations should perform acceptably, however, we have yet tosee a study that has adequately captured all true relations.
- 7
- 8 We have added a few sentences (blue is new) to the Discussion subsection 4.3, "ML Use for Knowledge
 9 Discovery" where we further urge for the use of TUURTs (Temporal, Unseen, Ungaged Region Tests)'
 10 (*revised lines 914-925*):
- 11

12 While it is understandable that not every ML-SWT paper aims to explain physical processes, the 13 SWT community should agree on a baseline of tests that all ML-SWT models undergo to assess 14 model robustness and transferability. Specifically, we urge use of TUURTs (Temporal, Unseen, 15 Ungaged Region Tests) for future ML-SWT models as a helpful step towards better modeling 16 practices, increased model transparency and robustness (Fig.1). As stated in figure 1, for TUURTs, 17 testing for "unseen" cases means testing only within the developmental dataset, whereas testing for 18 "ungaged" cases means testing for new sites that have no data and have not been previously seen by 19 the model at all. Due to the difficulty of conducting spatial tests compared to temporal tests, few ML-20 SWT studies have applied one or two of the tests, and rarely all three (Topp et al., 2023; Hani et al., 21 2023, Souassi et al., 2023). For example, Siegel et al. (2023), a non-ML SWT paper, tested for 22 ungaged regions and unseen data but did not perform a temporal test. To our knowledge, Philippus et al. (2024), appears to be the only published SWT-ML study that applied TUURTs with some success. 23 24 We further encourage modelers to shift towards more generalizable models, which are in theory, 25 more capable of performing well across diverse scenarios and datasets, and stand to become 26 increasingly important with the unpredictability of climate extremes.

27 28

29 2. Evaluation of Data Requirements: The manuscript does not extensively discuss the challenges that 30 ML ST modelers are facing with. Different ML models have varying data requirements, but the review 31 does not thoroughly discuss the data needs for each type of model. For example, ML models are 32 dependent on data. If we compare the availability of streamflow observation data availability versus the 33 SWT observation data, we realize there is a massive gap here, which impacts the studies and reduces the 34 SWT model performances. I suggest, while the authors encouraging the researchers and water institutes to 35 collect more data, they add their comments on this issue and discuss how researchers can reduce the 36 impact of this problem in their models. 37

AUTHOR RESPONSE: We agree with the referee that issues remain with data requirement limitations.
 We propose adding a new 'Discussion' subsection, titled 'ML Data Requirements vs. Availability' stating the following:

42 *new section 4.2 ML Data Requirements vs. Data Availability (*revised lines 881-906*).

43 44 While, in recent years, access to hydrologic data has improved (Miller et al., 2016; CUAHSI, 2024), 45 data remains scarce for many hydrologic applications including SWT research, particularly because 46 continual project management and funding to place and maintain stream temperature sensors, can be 47 expensive and/or time-consuming to undertake. As a result, in the 21st century, the scarcity of data 48 remains a large impediment for the application of machine learning in SWT modeling. What is more, 49 the question of data quantity (how much data do you have?) versus quality (how much diverse data is 50 needed?) continues to hinder ML use in hydrologic applications. Xu and Liang (2021) make the 51 excellent point that one year of streamflow data (can swap for stream temperature) at 15-minute 52 intervals equals about ~35,000 points, which may seem extensive, but is unlikely to be enough to 53 properly train a ML model due to autocorrelation and limited exposure to diverse types of data that

are naturally encountered with a longer time-series (Xu and Liang, 2021). For example, machine
learning models may only predict flood volumes they have previously seen (Kratzert et al., 2019).
While data requirements for ML remain high, there are some strategies that researchers have used to
alleviate this impact.

5 One strategy that hydrologists in other fields have used to tackle this problem is data 6 augmentation, which can be applied spatially or temporally to create new training examples that the 7 ML model can learn from. Spatial augmentation can be done by means of interpolation methods, i.e., 8 kriging or distance weighting to create new data points or by generating synthetic data based on 9 expected physical patterns to fill gaps in data coverage (Baydaroğlu and Demir, 2024). Temporal data 10 augmentation can be done by shifting, scaling or adding noise to existing time series to create new training examples (Skoulikaris et al., 2022). Alternatively, and not a new idea, is to use the statistical 11 12 technique known as seasonal decomposition, which breaks down a time series into its main 13 components, i.e., the trend, seasonal patterns and residual components (Apavdin et al., 2021; He et 14 al., 2022). These can then be recombined to generate new data and train the model for improved 15 accuracy (Apaydin et al., 2021). In addition to data augmentations, data requirements can be 16 alleviated by considering the help of unsupervised transfer learning, i.e., use pre-trained models on 17 similar tasks to reduce amount of data needed for training, or semi-supervised learning, such as few 18 shot learning, i.e., combine a small percent of labeled data with larger percent of unlabeled data to 19 improve model performance (Yang et al., 2023). By implementing these strategies, researchers in 20 other hydrologic fields have shown that models can be improved with less data, strategies that are 21 likely transferable to SWT research.

22 23

24 3. Future Directions Could Be Expanded: Although the paper concludes with a general discussion of 25 future challenges, it does not offer specific, actionable directions for future research. Highlighting key 26 areas where ML can advance, such as the use of satellite data, sensor networks, or the fusion of climate 27 models with ML, would provide more meaningful insights. In this concept, we can learn from hydrologic 28 community and capitalize on their experience and what they learned. The ML hydrologic community is 29 moving toward making global models, incorporating mechanistic models into their ML framework and 30 learning the governing factors, flow prediction with predicted inputs (predicted meteorological inputs) 31 and last but not least, providing a seamless simulation in streams in CONUS/global scale. Therefore, I 32 would ask the authors to add their comments on where the future direction of SWT community should be 33 and how SWT community can achieve the future objectives and what the barriers are. 34

AUTHOR RESPONSE: We agree and appreciate the referee's feedback. We propose adding a new
 'Discussion' subsection, titled '<u>4.4 Future Directions of SWT Modeling</u>', with the following:

*new Section 4.4 Future Directions of SWT Modeling, (*revised lines 944-991*):

40 The utility of ML in hydrologic modeling has advanced significantly, with interest seemingly 41 growing exponentially (Nearing et al., 2021). With the novelty of ML, it is easy to over-value model 42 performance and ignore the physics of the system, but with several decades of ML-experience, we 43 advocate it is necessary to purposefully use ML to address physically-meaningful questions and not 44 just create ML for the sake of creating. Given this, Varadharajan et al. (2022) laid out an excellent 45 discussion on opportunities for advancement of ML in water quality modeling, see section 3 of 46 publication of Varadharajan et al. (2022). (Varadharajan et al., 2022)Here we highlight some of the 47 questions from Varadharajan et al. (2022) that can be considered in the context of what objectives the 48 SWT community should be using in the ML era, namely: 1) How do we use physical knowledge (re: 49 heat exchange equations, radiation influence) to improve models and process understanding? 50 Rahmani et al. (2023) coupled NNs with the physical knowledge from SNTEMP, a one-dimensional 51 stream temperature model that calculates the transfer of energy to or from a stream segment by either 52 heat flux equations or advection, but found that even with SNTEMP, their flexible NNs exhibited 53 substantial variance in prediction and needed to be constrained by further multi-dimensional

 assessments (Rahmani et al., 2023). In short, if our use of physics in machine learning makes our models worse, we should understand why.

3 A second question that needs addressing is 2) How do we deal with predictive uncertainty in ML 4 used for SWT modeling? According to Moriasi et al. (2007), uncertainty analysis is the process of 5 quantifying the level of confidence in any given model output based on five guidelines: 1) the quality 6 and amount of observations (data), 2) the lack of observations due to poor or limited field monitoring, 7 3) the lack of knowledge of physical processes or operational procedures (instrumentation), 4) the 8 approximation of our mathematical equations, and 5) the robustness of model sensitivity analysis and 9 calibration. For example, in rainfall-runoff modeling, researchers have proposed benchmarking to 10 examine uncertainty predictions of ML rainfall-runoff modeling (Klotz et al., 2022). For stream temperature modeling, researchers have attempted to address the role of uncertainty in deep learning 11 12 model (RGCN, LSTM) predictions using the Monte Carlo Dropout (Zwart, Oliver, et al., 2023) and a 13 unimodal mixture density network approach (Zwart, Diaz, et al., 2023).

14 Other questions that SWT-ML studies should consider is 3) How do we make ML models 15 generalize better, specifically with regards to ungaged basins? And 4) How can ML models be 16 improved to predict extremes? As ML models advance to use satellite data, include more sensor 17 networks and/or couple with climate models, there is a logical next step toward creating generalizable 18 models that can account for extremes. In our review, only two papers by the same group (Rahmani et 19 al., 2020, 2023) conducted a CONUS-scale approach towards SWT-ML modeling, omitting 20 hydrologically important regions in the southwest (CA) and southeast (FL). Recently, a satellite 21 remote sensing paper used RF to model monthly stream temperature across the CONUS and tested for 22 temporal (walk-forward validation), unseen and 'true' ungaged regions (Philippus et al., 2024). We 23 have also learned that ML models such as LSTMs, generally only make predictions within the bounds 24 of their training data (Kratzert et al., 2019), which is a limitation for predicting extremes. Thus, we 25 strongly urge the community to work towards ML models that generalize better and/or are more 26 robust towards predictions of extremes.

27 Finally, 5) How can we build ML models such that they are seen as trustworthy and 28 interpretable by the hydrologic community? To answer this question, we must address a technical 29 barrier (black-box issues, data limitations, model uncertainty) and a social barrier (i.e., educated 30 skepticism of ML due to novelty, little understanding of computer science basics and/or coding 31 experience). If we are to incorporate ML into decision-making processes, it makes sense that ML 32 must be transparent and understandable to more than just computer or data scientists (Varadharajan et 33 al., 2022). For example, Topp et al. (2023) recently used explainable AI to elucidate how ML 34 architectures affected the SWT model's spatial and temporal dependencies, and how that in turn 35 affected the model's accuracy. Addressing this technical barrier can also be done by improving access 36 to data, which has seen remarkable progress thanks to web repositories such as NSF-funded 37 CUAHSI's Hydro share (CUAHSI, 2024) and GitHub (GitHub, 2024). In the United States, data 38 access to state and locally-based data remains limited, and should be addressed. In terms of the social 39 barrier, education about ML and ML-use is key. Societal interest in ML has thankfully also lead to a 40 plethora of educational resources and ML walk-through videos and tutorials in Tensorflow (Abadi et 41 al., 2016), PyTorch (Paszke et al., 2019), and Google Colab (Bisong, 2019). With the speed at which 42 ML-use is evolving, short communication pieces (Lapuschkin et al., 2019) and opinion pieces 43 (Kratzert et al., 2024) with clear examples about an ML-issue and practical solutions will also help 44 make ML challenges more transparent and therefore accessible to the hydrologic community-at-45 large.

46

47 <u>Added citations used for new subsection, 4.4 Future Directions of SWT Modeling:</u>

48 1) Apaydin, H., Taghi Sattari, M., Falsafian, K., and Prasad, R.: Artificial intelligence modelling integrated with
 49 Singular Spectral analysis and Seasonal-Trend decomposition using Loess approaches for streamflow
 50 predictions, Journal of Hydrology, 600, 126506, https://doi.org/10.1016/j.jhydrol.2021.126506, 2021.

2) Baydaroğlu, Ö. and Demir, I.: Temporal and spatial satellite data augmentation for deep learning-based rainfall nowcasting, Journal of Hydroinformatics, 26, 589–607, https://doi.org/10.2166/hydro.2024.235, 2024.

- 3) CUAHSI. 2024. Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI) Water
 Data Portal: https://www.cuahsi.org/community/water-data-portals, last access: 13 November 2024.
- 4) Kratzert, F., Gauch, M., Klotz, D. and Nearing, G., 2024. HESS Opinions: Never train an LSTM on a single basin. Hydrology and Earth System Sciences Discussions, 2024, pp.1-19.
- 5) Kwak, J., St-Hilaire, A., and Chebana, F.: A comparative study for water temperature modelling in a small basin,
 the Fourchue River, Quebec, Canada, Hydrological Sciences Journal, 1–12,
 https://doi.org/10.1080/02626667.2016.1174334, 2016.
- 8 6) Philippus, D., Sytsma, A., Rust, A., and Hogue, T. S.: A machine learning model for estimating the temperature of small rivers using satellite-based spatial data, Remote Sensing of Environment, 311, 114271, https://doi.org/10.1016/j.rse.2024.114271, 2024.
- 7) Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C., and Gupta, H. V.:
 What Role Does Hydrological Science Play in the Age of Machine Learning?, Water Resources Research, 57, e2020WR028091, https://doi.org/10.1029/2020WR028091, 2021.
- 8) Skoulikaris, C., Venetsanou, P., Lazoglou, G., Anagnostopoulou, C., and Voudouris, K.: Spatio-Temporal Interpolation and Bias Correction Ordering Analysis for Hydrological Simulations: An Assessment on a Mountainous River Basin, Water, 14, 660, https://doi.org/10.3390/w14040660, 2022.
- 9) Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: A Simple Way to
 Prevent Neural Networks from Overfitting, Journal of Machine Learning Research, 15, 30, 2014.
- 10) Yang, M., Yang, Q., Shao, J., Wang, G., and Zhang, W.: A new few-shot learning model for runoff prediction:
 Demonstration in two data scarce regions, Environmental Modelling & Software, 162, 105659, https://doi.org/10.1016/j.envsoft.2023.105659, 2023.
- 11) GitHub. 2024. About Git and Github: https://docs.github.com/en/get-started/start-your-journey/about-github and-git, last access: 14 November 2024.
- 12) Lapuschkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W. and Müller, K.R., 2019. Unmasking Clever
 Hans predictors and assessing what machines really learn. Nature communications, 10(1), p.1096.
- 26 13) Zwart, J.A., Oliver, S.K., Watkins, W.D., Sadler, J.M., Appling, A.P., Corson-Dosch, H.R., Jia, X., Kumar, V. and
 27 Read, J.S., 2023. Near-term forecasts of stream temperature using deep learning and data assimilation in support
 28 of management decisions. JAWRA Journal of the American Water Resources Association, 59(2), pp.317-337.
- 29 14) Zwart, J.A., Diaz, J., Hamshaw, S., Oliver, S., Ross, J.C., Sleckman, M., Appling, A.P., Corson-Dosch, H., Jia,
 30 X., Read, J. and Sadler, J., 2023. Evaluating deep learning architecture and data assimilation for improving
 31 water temperature forecasts at unmonitored locations. *Frontiers in Water*, 5, p.1184992.
- 15) Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S. and
 Nearing, G., 2022. Uncertainty estimation with deep learning for rainfall-runoff modeling. Hydrology and
 Earth System Sciences, 26(6), pp.1673-1693.
- 35 16) M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S.
 36 Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, R. Jozefowicz, Y. Jia, L. Kaiser, M. Kudlur, J.
 37 Levenberg, D. Mané, M. Schuster, R. Monga, S. Moore, D. Murray, C. Olah, J. Shlens, B. Steiner, I. Sutskever,
 38 K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke,
 39 Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems. 2015. TensorFlow.
 40 Website: https://www.tensorflow.org/
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. Website: https://arxiv.org/abs/1912.01703
- 45 18) Bisong, E. (2019). Google Colaboratory. In: Building Machine Learning and Deep Learning Models on Google
 46 Cloud Platform. Apress, Berkeley, CA. Website: <u>https://doi.org/10.1007/978-1-4842-4470-8_7</u>
- 47

1 4. The manuscript walked through many ML and AI models. An important factor of the ML and AI

2 models are the inputs. I assume you faced a variety of inputs that have been used in the models. That

3 would be informative to the readers, if the authors add their observations that what kind of inputs that

4 have been missed to be used, either because it is not available yet or it is even missed. For instance,

whether there is any geophysical attribute, climatic attributes, or any forcings that is worth to be extractedand used in ML models.

AUTHOR RESPONSE: We appreciate the referee's feedback. In the Supplementary Materials, Table S1
 contains some of the suggested data by the referee, such as: period considered, region examined, temporal
 resolution of SWT, spatial scale of study, and hydrometeorological parameters used for modeling. We
 provided the information as Supplementary Material because Tables S1 and S2 are seven pages alone,
 which may risk making the review lengthier than it already is. We have added text to the manuscript
 regarding model inputs and moved the LASSO paragraph (original lines 247-253) to this section because

13 14

15 *new section 2.4.2 Model Inputs for ML-SWT (*revised lines* 488-516):

we think it can more smoothly follow the paragraph on feature importance.

16

17 Using air temperature (AT) to better understand SWT has been considered since the 1960s, when 18 Ward (1963) and Edinger et al. (1968) discussed the influence of air temperature on SWT. Since then, 19 studies have used varying input variables (see Table S1), however, the model inputs of AT and SWT 20 continue to be the most used in ML-modeling studies. In particular, studies have used AT from time 21 periods outside of the known SWT record to improve model performance (Sahoo et al., 2009; 22 Piotrowski et al., 2015; Graf et al., 2019). In addition to AT and SWT, flow discharge has been used 23 to attempt to constrain SWT (Foreman et al., 2001; Tao et al., 2008; St-Hilaire et al., 2011; Grbić et al., 2013; Piotrowski et al., 2015; Graf et al., 2019; Qiu et al., 2020). Traditionally-used model inputs 24 25 include precipitation (Cole et al., 2014; Jeong et al., 2016; Rozos, 2023), wind direction/speed (Hong 26 and Bhamidimarri, 2012; Cole et al., 2014; Jeong et al., 2016; Kwak et al., 2016; Temizyurek and 27 Dadaser-Celik, 2018; Abdi et al., 2021; Jiang et al., 2022), barometric pressure (Cole et al., 2014), 28 landform attributes (Risley et al., 2003; DeWeber and Wagner, 2014; Topp et al., 2023; Souaissi et 29 al., 2023), and many more (see Table S1).

30 In the last few years, including the day-of-year as an input, DOY (Oiu et al., 2020; Heddam et 31 al., 2022; Drainas et al., 2023; Rahmani et al., 2023) and humidity (Cole et al., 2014; Hong and 32 Bhamidimarri, 2012; Kwak et al., 2016; Temizyurek and Dadaser-Celik, 2018; Abdi et al., 2021), 33 have also shown to better capture the seasonal patterns of SWT (Qiu et al., 2020; Philippus et al., 34 2024). With improved access to remote sensing data, there has also been a notable increase of satellite 35 products such as estimates of sky cover (Cole et al., 2014), solar radiation (Kwak et al., 2016; Topp et 36 al., 2023; Majerska et al., 2024), sunshine per day (Drainas et al., 2023) and potential ET (Rozos, 37 2023; Topp et al., 2023). However, more research is needed to better understand the influence of 38 newer model inputs on SWT (Zhu and Piotrowski, 2020).

39 Recently, SWT studies focused on the CONUS-scale have chosen to use as many model inputs 40 as available, with Wade et al. (2023), a point-scale CONUS ML study using over 20 variables, while 41 Rahmani et al. (2023) created a LSTM model and considered over 30 variables to simulate SWT. 42 Despite the use of diverse data, the models in these studies performed only satisfactorily and were 43 deemed not generalizable, leaving much room for improvement in CONUS-scale modeling of SWT. 44 With the compilation of larger and larger datasets, feature importance in ML, that is the process of 45 using techniques to assign a score to model input features based on how good the features are at 46 predicting a target variable, can be an efficient way to improve data comprehension, model 47 performance, and model interpretability, the latter of which can dually serve as a transparency marker 48 of which features are driving predictions. Methods for measuring feature importance include using 49 correlation criteria (Pearson's r, Spearman's rho), permutation feature importance (shuffling feature 50 values, measuring decrease in model performance), linear regression feature importance (larger absolute values indicate greater importance), or if using CART/RF/gradient boosting, entropy 51 52 impurity measurements can be insightful (Venkateswarlu and Anmala, 2023). 53

Moved part of section 2.3.1, original (lines 246-253) to section 2.4.2 Model Inputs for ML-SWT (moved to lines 517-523):
 For example, one technique that can be used to improve ML model parameter selection is the

For example, one technique that can be used to improve ML model parameter selection is the *Least Absolute Shrinkage and Selection Operator (LASSO)*, a regression technique used for feature selection (Tibshirani, 1996). Research utilizing ML models for SWT frequency analysis at ungaged basins used the LASSO method to select explanatory variables for two ML models (Souaissi et al., 2023). The LASSO method consists of a shrinkage process where the method penalizes coefficients of regression variables by minimizing them to zero (Tibshirani, 1996). The number of coefficients set to zero depends on the adjustment parameter, which controls the severity of the penalty. Thus, the method can perform both feature selection and parameter estimation, an advantage when examining large datasets (Xu & Liang, 2021).

- 13 5. Lack of Clear Structure in the Evaluation: Although the paper aims to summarize the performance 14 evaluation metrics for ML models in SWT prediction, the organization of these sections feels somewhat 15 scattered. A more systematic approach could improve clarity, such as separating the analysis based on 16 time scales (e.g., hourly, daily, monthly) or spatial scales (local, regional, continental). This would make 17 it easier for readers to find the relevant insights based on their application. For instance, a stream 18 temperature model in monthly scale is different from a daily or hourly scale models on many aspects. As 19 an example, the complexity of a daily model is different from a monthly temperature models. A monthly 20 model may not need all inputs of a daily model to capture the monthly changes. The authors can add their 21 overall opinion of what types of models are better fitted to which time scale. In ML models, it is 22 important to know the scope of the model, whether it is a local model that needs to be calibrated site by 23 site, or it is a model that is designed to work for multiple sites (a regional model). I believe that would be
- 24 informative to consider the modeling approach when methods are compared.
- 25

4

5

6

7

8

9

10

11

12

AUTHOR RESPONSE: We appreciate the opportunity to clarify. Initially, we compiled a performance
 metric comparison by spatial scale for the most-cited metric, RMSE (42 papers cited) and plotted RMSE
 by study for regional/CONUS scale and local scale (located in HYDROSHARE repository but not the
 manuscript, see plots on next page for comparison). The comparison found minimal difference in RMSE
 between the regional/CONUS studies and the local scale studies, which we summarize in Table 1.

31

32 Given the speed at which ML is advancing and being applied for hydrologic applications, we do not think 33 it wise to opinionate on which ML model is better or worse. Instead, our aim is to inform the reader on

34 current studies and metrics, which we do by providing data as supplementary info, such as Tables S1

35 which states the time scale, spatial scale, region and time period considered of each study while Table S2

36 lists the data analysis techniques and/or ML algorithms used, as well as the training/validation/testing

37 percentages/time periods as reported by the study. We think that summarizing publications (see Tables S1

and S2) and compiling performance metrics allows the reader to identify what has already been done in

39 the ML-SWT field so that they can then make their own informed decisions depending on their research

40 questions, model selection, project time frame, etc.



RMSE, Local-scale



- Foreman_etal2000_Fraser_riv_Canada
 Hong_Bhadimirri_2012_NewZealand
- Hong_Bradimirri_2012_NewZealand Herbert_etal2014_Oregon
- Cole_etal2014_Delaware
- Jeong_etal2016_Soyang_Riv_SKorea
- Laanaya_etal2017_StMarguerite_Canada
 Zhu_Heddam_Wu_etal2019_WA_OR_NY
- Zhu_Heddam_Nyarko_etal_2019_Europe
- Lu_Ma2020_Tualatin_riv_Oregon
- Khosravi_etal2023_Central_Delaware_riv
- Zhu_Heddam_2019

- StHilaire_etal2011_Moisie_riv_Canada
- Grbic_etal2013_Drava_Croatia
- HadzimaNyarko_etal2014_Drava_Croatia
- 🔲 Rabi_etal2015_Drava_Croatia
- Kwak_etal2016_Fourchue_Canada
- TemizyurekDadaserCelik2018_Kizilm_Turkey
- Zhu_Nyarko_HadzimaNyarko_2019
- Zhu_Nyarko_Gao_etal_2019_Drava
- Zanoni_etal2022_Italy
- Wade_etal2023_two_crks_Oregon

1 6. The authors need to decide first who are the readers of the papers. Whether the paper serves to new-

- commers to ML and AI methodologies in stream temperature community or it serves to researchers that
 are already familiar with basics of ML and AI methods.
- 4

5 **AUTHOR RESPONSE:** We agree with the referee that the purpose of the review should be more clearly

stated. We drafted this paper to serve as a middle ground between traditional modelers and more wellversed ML users. The intended audience are hydrologic modelers who have heard of AI/ML and want a

- 8 summary of what has been done in SWT modeling using ML. Our dual objective is also for this to be a
- 9 reference for assessment of ML performance. At the same time, we want ML researchers to be aware of
- 10 where their models stand compared to other modelers while communicating that an "A+ grade" is
- 11 actually more common (and therefore the new average) relative to what they are used to in hydrologic
- 12 modeling. We have added a few sentences in the introduction, under section 1.2 'Study Objectives' of the 13 manuscript to state who the intended audience is:
- 14 **1.2. Study Objective** (new in blue, *revised lines* 64-68)

15 The current work includes an extensive literature review of studies that used ML algorithms/models for 16 river/SWT modeling, hindcasting and forecasting. The intent of this review is two-fold: 1) to introduce ML for hydrologists who have modeling experience and are interested in pursuing ML-use for their SWT 17 18 studies, and 2) to provide a broad overview of machine learning applications in SWT. For ML experts, 19 we think that this review could also prove useful as reference for how ML has been applied in the field 20 of SWT modeling and where improvement is needed. Overall, this article aims to serve as a bridge 21 between hydrologists and machine learning experts. Our review includes papers cited by Zhu and 22 Piotrowski (2020), who previously conducted a study of ANNs used in SWT modeling, however, we 23 provide a comprehensive examination of peer-reviewed journals that use any type of artificial 24 intelligence/ML algorithm to model or evaluate river/SWT [...]

25 26

7a. While the paper provides an extensive review of ML applications in SWT modeling, it focuses
heavily on listing the types of ML models used rather than deeply analyzing their applications, strengths,
weaknesses, and performance differences. A more critical analysis of the pros & cons of each model type
could provide greater value to researchers choosing the appropriate model for their specific needs. To
provide a few examples, I refer to lines 136 – 143 & lines 146 – 159 & lines 263 - 292.

33 **AUTHOR RESPONSE:** Thank you for the opportunity to clarify. We provided supplementary tables to 34 summarize study information, for example, Tables S1 includes summarized information stating the time 35 scale, spatial scale, region and time period considered of each study while Table S2 lists the data analysis techniques and/or ML algorithms used, as well as the training/validation/testing percentages/time periods 36 as reported by the study. We think the "pros/cons" and "strengths/weakness" vary depending on the 37 38 research goal and question, and the robustness of ML models allows them to cater to most problems, 39 which is why rather than opinionating, we provide concrete specifications on the models used and allow 40 the reader to decide based on their objectives.

- 41
- 42

7b. The first half of the paragraph that is written in lines 136 – 143 explains the fundamentals of the
method, which may not be necessary to be long, and the rest is an example of the method usage.
However, this paragraph could have been enriched by statements like the advantages and disadvantages
of this method compared to other existing ML methods or even to a linear regression method, or a 1D
mechanistic method (although they are not ML methods, but the comparison is beneficial to the readers).
The authors also can add their statement of under what conditions they think the method is beneficial.

AUTHOR RESPONSE: We agree with the reviewer. We have revised the text to describe the
 advantages and disadvantages of K-nn, section 2.3.1.1 (crossed out is deleted, *revised lines 166-181*):

1	
2	K-nearest neighbors (K-nn) is a type of versatile supervised ML algorithm (Fix & Hodges, 1952;
3	Cover & Hart, 1967) used to solve nonparametric classification and regression problems. It is one of
4	the oldest algorithms (Fix & Hodges, 1952; Cover & Hart, 1967) considered within classical ML. The
5	K-nn algorithm uses proximity between data points to make classifications or evaluations about the
6	grouping of any given data point (Acito, 2023). K-nn gained popularity in the 2000s due to its
7	simplicity in implementation and understanding, making it readily accessible to hydrologic
8	researchers and practitioners. While less used today, For example, StHilaire et al. (20122011) used
9	various K-nn model configurations to model SWT for the Moisie River in northern Quebec, Canada,
10	finding that. T the best K-nn model required prior-day SWT data and day-of-year (DOY), an indicator
11	of seasonality (St. Hilaire et al., 2011). Advantages of K-nn include its non-assumptions of the
12	underlying distribution of the data, allowing it to handle nonlinear complexities without requiring a
13	solid model structure as is the case for some physical models (St. Hilaire et al., 2011). Disadvantages
14	of K-nn are that it is computationally intensive, may require extensive cross-validation, performance
15	can be affected by irrelevant/redundant features, and due to its high memory and computational
16	needs, is impractical for large-scale applications, i.e., scalability issues, (Acito, 2023). For example,
17	Heddam et al. (2022) For example, Heddam et al. (2022) For five stream stations in Poland, Heddam
18	et al. (2022) compared K-nn with other ML algorithms, finding that K-nn was outperformed by other
19	MLs such as least squares support vector machine and neural networks. performed poorly compared-
20	to other ML algorithms. The use of K-nn may still be reasonable for simple, local cases but we advise
21	considering other MLs for more complex or larger-use cases.
22	

- 7c. Lines 146 153 explains PCA & ck-means clustering on data reduction application, however, it is not
 clear under what conditions we can use them.
 - AUTHOR RESPONSE: We agree. We propose adding text to clarify:
- **29** *Added lines 192-193:*

Krishnaraj and Deka (2020) used *K-means* to organize spatial grouping for water quality monitoring stations for dry and wet regions along the Gangas River basin in India to identify whether pollution patterns could be discerned.

34 *Added lines 198-203:*

35 Using PCA, Krishnaraj and Deka (2020) found that certain water quality parameters (dissolved 36 oxygen, sulfate, electrical conductivity) were more dominant in the dry season compared to the wet 37 season (total dissolved solids, sodium, potassium, sodium, chlorine, chemical oxygen demand), data 38 which could be used to cater the monitoring program to the important parameters. In their study, 39 SWT was not a dominant parameter, likely in part because the SWT of large downstream rivers like 40 the Gangas River are generally less variable due to their larger volume and stronger thermal buffer. 41 Used k-means and PCA in the Ganga River Basin of India to find spatiotemporal patterns of water 42 quality parameters, including SWT.

43 44

53

23

27

28

30

31

32

- 45 7d. Additionally, that would be nice for readers if the authors add feature importance to their comparison
 46 as it has been used more frequently in streamflow and soil moisture prediction studies.
 47
- 48 AUTHOR RESPONSE: We agree and added text on feature importance to a section on model inputs as
 49 suggested (please see comment #4 for full text).
 50
- 51 The text specific to feature importance is below (*revised lines 506-516*):52
 - Recently, SWT studies focused on the CONUS-scale have chosen to use as many model inputs

1 as available, with Wade et al. (2023), a point-scale CONUS ML study using over 20 variables, while 2 Rahmani et al. (2023) created a LSTM model and considered over 30 variables to simulate SWT. 3 Despite the use of diverse data, the models in these studies performed only satisfactorily and were 4 deemed not generalizable, leaving much room for improvement in CONUS-scale modeling of SWT. 5 With the compilation of larger and larger datasets, feature importance in ML, that is the process of 6 using techniques to assign a score to model input features based on how good the features are at 7 predicting a target variable, can be an efficient way to improve data comprehension, model 8 performance, and model interpretability, the latter of which can dually serve as a transparency marker 9 of which features are driving predictions. Methods for measuring feature importance include using 10 correlation criteria (Pearson's r, Spearman's rho), permutation feature importance (shuffling feature values, measuring decrease in model performance), linear regression feature importance (larger 11 12 absolute values indicate greater importance), or if using CART/RF/gradient boosting, entropy 13 impurity measurements can be insightful (Venkateswarlu and Anmala, 2023).

14 15

28

29

30

34 35

36

37 38

39

40 41

7e. Lines 263 – 292 are organized in three paragraphs while providing general knowledge about ANNs
 with relatively less direct relations to water temperature application.

AUTHOR RESPONSE: We appreciate the reviewer's feedback and are open to making changes to improve the manuscript for the reader. Referee #3 made a similar comment about this section, and we now wonder if it would be better to provide the description of ANN variants and alternatives (lines 263-320) as part of an appendix. We think it would still be helpful to keep the information, but we also agree that it may be too extensive for the main text. In this way, the manuscript can be made more concise while also keeping the details as a section of the manuscript for anyone who is interested in reading further.

Following this line of thinking, we can point the reader to the appendix (*revised lines 326-327*):

"For more detail on traditional ANNs, with descriptions of ANN variants and backpropagation alternatives, we refer the reader to Appendix A."

We have added Appendix A after the conclusion (*revised lines 1014 – 1075*).

33 Minor corrections:

1. Line 13: There is a typo that changes the meaning of the sentence. It should be "... with in situ ..." or "... with in-situ ...".

AUTHOR RESPONSE: Thank you for pointing this out, we have fixed the typo to read "with in-situ" (*revised line 13*).

42 2. Line 132: There is a typo here too. It should be "long short-term memory". Although I am trying to catch them, there is a chance that I miss some of them. I recommend the authors to carefully re-read the manuscript or ask help from a fresh pair of eyes to find these types of typos.
45

46 AUTHOR RESPONSE: Thank you! We have revised the text to read "long short-term memory"
 47 (*revised line 160*) and reviewed the text accordingly.

48 49

50 3. Lines 208 – 210: to make the sentence more accurate, it needs to be stated whether these are local models or one model for multiple sites. Additionally, I believe by "NNs" here, the authors mean feedforward neural network, which are totally different from recurrent neural networks.
53

AUTHOR RESPONSE: Yes, we agree with both points. We have clarified that a feed-forward NN was 3 used and revised the sentence to make it more accurate (revised lines 263-266): In the case of A SWT modeling study comparing the output of three model versions of DT, GPR, and feed-forward neural networks for daily SWT modeling multiple sites and prediction, found that DTs can could perform similarly to GPR and feed-forward neural networks when detailed statistics of air temperature, day-of-year, and discharge were included NNs (Zhu, Nyarko, Hadzima-Nyarko, Heddam, et al., 2019). 4. Line 541: "at" is missed. It is .. All journals examined used at least ..." AUTHOR RESPONSE: Thank you! We have added the word "at" (revised line 681).

1 <u>Referee #2 Comments</u>

- 2 This is a meaningful manuscript that provides a thorough review of ML approaches for SWT modeling
- 3 and their evaluation metrics. I believe that the current scientific community has indeed developed a broad
- 4 understanding of the integration of ML into stream temperature modeling. Hence, while the manuscript
- 5 presents a comprehensive overview, incorporating more in-depth insights could enhance its appeal to
- 6 readers and significantly increase its contribution to the field. The review covers a wealth of content,
- 7 including recent articles and other reviews, but the sections are somewhat loosely structured, with key
- 8 points relatively briefly mentioned.

9 AUTHOR RESPONSE: We thank the referee for their time and feedback, we believe the manuscript is

- stronger as a result. We address specific referee comments below. For reference, we separated some referee comments into a, b, etc., to provide a more organized response. Proposed new/edited text is in
- 12 BLUE. Revised lines in the track-changes manuscript are indicated by: (*revised lines XXX-XXX*).
- 13 14

15 1. For instance, in the first section (Overview: SWT Model Types), the author provides a solid overview

- 16 of statistical, physical, and ML models. However, a more detailed analysis of the comparative strengths
- and weaknesses of physical and ML models would strengthen the discussion. The models are presented in
- 18 a nearly linear developmental order in this review, but it would be beneficial to mention some points, for
- 19 example, [if] physical models perform well, why ML models are adopted[?].
- AUTHOR RESPONSE: The referee makes a good point with regards to the question of "if physical
 models perform well, why are ML models being adopted?". We have expanded the section 2.3 "Artificial
 Intelligence Models in SWT Modeling" to discuss this (*revised lines 136-156*):

22

24 In the last decade, computing advances in AI have started to offer several advantages for using machine 25 learning (ML) in hydrology that are comparable to physically based models (Cole et al., 2014; Zhu et 26 al., 2019; Rehana and Rajesh, 2023). In contrast to traditional physically based models, the code 27 underlying ML models are generally open-source and publicly available allowing for near real-time 28 accessible advances and user feedback, whereas the source code for some physically based models may 29 be inaccessible to the public due to being privately managed (MIKE suite of models) or the model 30 software may be publicly available but take years to publish updates (USGS MODFLOW, Simunek's 31 HYDRUS). One advantage that has made ML increasingly appealing includes its ability to learn 32 directly from the data (i.e., data driven), which can be useful when the underlying physics are not fully 33 understood or are considered too complex to model accurately.

34 Additionally, ML models are more efficient in making predictions compared to the time-intensive 35 solvers of physically based models. ML models can also handle the challenge of scalability, that is 36 managing large datasets and seamlessly deploying across various computer platforms and applications 37 (Rehana and Rajesh, 2023). Air2stream, a hybrid statistical-physically based SWT model (Toffolon and 38 Piccolroaz, 2015; Piccolroaz et al., 2016), initially outperformed earlier ML models such as Gaussian Process Regression (Zhu et al., 2019). However, in the last few years, Air2stream has had its 39 40 performance matched and even exceeded by recent neural networks models (Feigl et al., 2021; Rehana 41 and Rajesh, 2023).

42 Finally, with computer processing power improving and the emergent field of quantum computing, 43 there is a strong belief that using ML and by extension AI, in science applications will drive innovation 44 to the point where natural patterns and insights not currently apparent in physical modeling will be 45 uncovered (Varadharajan et al., 2022). Thus, while physically based models are considered tried-and-46 true, thereby invaluable for their interpretability and grounding in established physics, ML models have 47 the potential for growth – where they can be used to first complement and eventually lead as powerful 48 tools for prediction, optimization, and understanding in increasingly complex and data-rich 49 environments. 50

51 New citation:

Toffolon, M. and Piccolroaz, S., 2015. A hybrid model for river water temperature as a function of air temperature and discharge. *Environmental Research Letters*, *10*(11), p.114011.

2 3

1

2. How to gain the trust of traditional model users in ML methods? (This question is inherently
challenging, as model users often have preferences based on their own familiarity with certain models and
may exhibit biases against alternative approaches. However, it may be worthy to acknowledge this in the
review.) This discussion could extend to the choice between different ML models as well, as conclusions
favoring one model over another often depend on the specific context of the study. Many conclusions are
applicable only under particular circumstances, so a generalization such as "a certain model is better

- 10 suited to a particular type of problem" is more appropriate.
- AUTHOR RESPONSE: We agree and appreciate the referee's feedback. We address this comment in our response to referee #1 for comment #3 (copied below) where we discuss how researchers can work to present their ML models as trustworthy. For this, we propose adding a new 'Discussion' subsection titled '4.4 Future Directions of SWT Modeling' (*revised lines 944-991*):
- 15

16 The utility of ML in hydrologic modeling has advanced significantly, with interest seemingly 17 growing exponentially (Nearing et al., 2021). With the novelty of ML, it is easy to over-value model 18 performance and ignore the physics of the system, but with several decades of ML-experience, we 19 advocate it is necessary to purposefully use ML to address physically-meaningful questions and not 20 create ML for the sake of creating. Given this, Varadharajan et al. (2022) laid out an excellent 21 discussion on opportunities for advancement of ML in water quality modeling, see section 3 of 22 publication Varadharajan et al., (2022). Here we highlight some of the questions from Varadharajan 23 et al. (2022) that can be considered in the context of what objectives the SWT community should be using in the ML era, namely: 1) How do we use physical knowledge (re: heat exchange equations, 24 25 radiation influence) to improve models and process understanding? Rahmani et al. (2023) coupled NNs with the physical knowledge from SNTEMP, a one-dimensional stream temperature model that 26 27 calculates the transfer of energy to or from a stream segment by either heat flux equations or 28 advection, but found that even with SNTEMP, their flexible NNs exhibited substantial variance in 29 prediction and needed to be constrained by further multi-dimensional assessments (Rahmani et al., 30 2023). In short, if our use of physics in machine learning makes our models worse, we should 31 understand why.

32 A second question that needs addressing is 2) How do we deal with predictive uncertainty in ML 33 used for SWT modeling? According to Moriasi et al. (2007), uncertainty analysis is the process of 34 quantifying the level of confidence in any given model output based on five guidelines: 1) the quality 35 and amount of observations (data), 2) the lack of observations due to poor or limited field monitoring, 36 3) the lack of knowledge of physical processes or operational procedures (instrumentation), 4) the 37 approximation of our mathematical equations, and 5) the robustness of model sensitivity analysis and 38 calibration. For example, in rainfall-runoff modeling, researchers have proposed benchmarking to 39 examine uncertainty predictions of ML rainfall-runoff modeling (Klotz et al., 2022). For stream 40 temperature modeling, researchers have attempted to address the role of uncertainty in deep learning 41 model (RGCN, LSTM) prediction using the Monte Carlo Dropout (Zwart, Oliver, et al., 2023) and a 42 unimodal mixture density network approach (Zwart, Diaz, et al., 2023).

43 Other questions that SWT-ML studies should consider is 3) How do we make ML models 44 generalize better, specifically with regards to ungaged basins? And 4) How can ML models be 45 improved to predict extremes? As ML models advance to use satellite data, include more sensor 46 networks and/or couple with climate models, there is a logical next step toward creating generalizable 47 models that can account for extremes. In our review, only two papers by the same group (Rahmani et 48 al., 2020, 2023) conducted a CONUS-scale approach towards SWT-ML modeling, omitting 49 hydrologically important regions in the southwest (CA) and southeast (FL). Recently, a satellite 50 remote sensing paper used RF to model monthly stream temperature across the CONUS and tested for temporal (walk-forward validation), unseen and 'true' ungaged regions (Philippus et al., 2024). We 51 52 have also learned that ML models such as LSTMs, generally only make predictions within the bounds 53 of their training data (Kratzert et al., 2019), which is a limitation for predicting extremes. Thus, we

strongly urge the community to work towards ML models that generalize better and/or are more
 robust towards predictions of extremes.

3 Finally, 5) How can we build ML models such that they are seen as trustworthy and interpretable 4 by the hydrologic community? To answer this question, we must address a technical barrier (black-5 box issues, data limitations, model uncertainty) and a social barrier (i.e., educated skepticism of ML 6 due to novelty, little understanding of computer science basics and/or coding experience). If we are to 7 incorporate ML into decision-making processes, it makes sense that ML must be transparent and 8 understandable to more than just computer or data scientists (Varadharajan et al., 2022). For example, 9 Topp et al. (2023) recently used explainable AI to elucidate how ML architectures affected the SWT 10 model's spatial and temporal dependencies, and how that in turn affected the model's accuracy. Addressing this technical barrier can also be done by improving access to data, which has seen 11 12 remarkable progress thanks to web repositories such as NSF-funded CUAHSI's Hydro share (CUAHSI, 2024) and GitHub (GitHub, 2024). In the United States, data access to state and locally-13 14 based data remains limited, and should be addressed. In terms of the social barrier, education about 15 ML and ML-use is key. Societal interest in ML has thankfully also lead to a plethora of educational 16 resources and ML walk-through videos and tutorials in Tensorflow (Abadi et al., 2016), PvTorch 17 (Paszke et al., 2019), and Google Colab (Bisong, 2019). With the speed at which ML-use is evolving, 18 short communication pieces (Lapuschkin et al., 2019) and opinion pieces (Kratzert et al., 2024) with 19 clear examples about an ML-issue and practical solutions will also help make ML challenges more 20 transparent and therefore accessible to the hydrologic community-at-large.

21

3a. Furthermore, the author may not clearly (separately) present the generalization capabilities of ML
 models in temporal and spatial contexts, which is crucial for data split. The model ability of

24 generalization over time is particularly meaningful for climate change studies, where overfitting (common

25 for ML studies) may lead to highly unreliable projections. Spatial generalization is useful for applying

26 models to new regions or watersheds (ungauged stream/river/watershed).

AUTHOR RESPONSE: We agree. Referee #1 made a similar comment (ref #1, comment #1A) about
overfitting and having ML undergo more testing and we propose to address both comments by adding: 1)
a new subsection 2.4.1, titled "Identifying model complexity", which discusses overfitting/underfitting,
with 2) a diagram with initial steps to mitigate overfitting. The new text is below:

32 *new Section 2.4.1, Identifying Model Complexity (*revised lines 464-483*)

33 34 The strong success of ML-use in SWT modeling warrants a brief and broad overview on identifying model complexity to minimize overfitting and underfitting" of models. When a model is too complex, 35 i.e., has too many features or parameters relative to the number of observations, or is forced to 36 37 overextend its capabilities, i.e., make predictions with insufficient training data, the model runs the 38 risk of overfitting (Srivastava et al., 2014). An overfitted model fits the training data "too well", 39 capturing noise and details that provide high accuracy on a training dataset, only to perform poorly 40 once the model encounters "unseen" data in testing/validation (Xu and Liang, 2021). Scenarios where 41 overfitting may be temporarily acceptable are: 1) model development is at preliminary stages, the 42 interest is in a "proof of life" concept, 2) when the objective is to identify heavily-relied on features 43 by the model, i.e., feature importance, or 3) in highly-controlled modeling environments where the 44 expected data will be consistently similar to the training dataset. The latter is more likely in industrial 45 applications and unlikely in the changing nature of hydrology. 46

In contrast, underfitting occurs when a model is too simple to capture any patterns in the data,
which can also lead to unsatisfactory performance in training, testing and validation. Underfitting can
occur with inadequate model features, poor model complexity or when regularization techniques,
(e.g., L1 or L2 regularization), are over-used, making the model too rigid and unable to respond to
changes in the data. Given the propensity of ML models to effectively learn the training data,
underfitting is less an issue in ML whereas overfitting can be widespread. In Figure 1, we present an

- 1 example workflow that researchers can use to transition away from overfitting and towards 2
 - generalizability. In the five-step outline (Fig. 1), we suggest the need for "Temporal, Unseen,
 - Ungaged Region Tests" (TUURTs), which is a call for temporal and spatially-focused testing that can
 - be used to strengthen model robustness.

Revised lines 484-486:



Figure 1. Diagram outlining steps that can be taken in modeling process to mitigate overfitting.

10 We propose to address the comment about generalization and a similar one made by ref #1 (comment #3) by adding a new Discussion subsection, '4.4 Future Directions of SWT Modeling'. Below is our selected 11

12 response where we state that models should work towards generalizability (revised lines 966-975). For 13 full text, please see comment #2:

14

7 8

9

15 Other questions that SWT-ML studies should consider are 3) How do we make ML models 16 generalize better, specifically with regards to ungaged basins? And 4) How can ML models be 17 improved to predict extremes? As ML models advance to use satellite data, include more sensor 18 networks and/or couple with climate models, there is a logical next step toward creating generalizable 19 models that can account for extremes. In our review, only two papers by the same group (Rahmani et 20 al., 2020, 2023) conducted a CONUS-scale approach towards SWT-ML modeling, omitting 21 hydrologically important regions in the southwest (CA) and southeast (FL). Recently, a satellite 22 remote sensing paper used RF to model monthly stream temperature across the CONUS and tested for temporal (walk-forward validation), unseen and 'true' ungaged regions (Philippus et al., 2024). We 23 24 have also learned that ML models such as LSTMs, generally only make predictions within the bounds 25 of their training data (Kratzert et al., 2019), which is a limitation for predicting extremes. Thus, we 26 strongly urge the community to work towards ML models that generalize better and/or are more

robust towards predictions of extremes.

3 3b. Additionally, the review does not systematically address the critical issue of model input selection,
which is essential in ML modeling. Model inputs for SWT modeling may include hydrometeorological
and physical parameters (or other attributes used in different studies), they play a role in model
performance and should be discussed in this part.

6 performance and should be discussed in this part.

AUTHOR RESPONSE: Thank you for pointing out this area in need of clarity. Referee #1, comment #4
had a similar question about model input, and we propose adding the paragraph below in response to
both. Additionally, we want to note that we included in Supplementary Materials, Table S1, which
contains some of the suggested data by the referee, such as: period considered, region examined, temporal
resolution of SWT, spatial scale of study, and hydrometeorological parameters used for modeling.

11 12

14

13 *new subsection 2.4.2, Model Inputs for ML-SWT (*revised lines 488-516*):

15 Using air temperature (AT) to better understand SWT has been considered since the 1960s, when 16 Ward (1963) and Edinger et al. (1968) discussed the influence of air temperature on SWT. Since then, 17 studies have used varying input variables (see Table S1), however, the model inputs of AT and SWT 18 continue to be the most used in ML-modeling studies. In particular, studies have used AT from time 19 periods outside of the known SWT record to improve model performance (Sahoo et al., 2009; 20 Piotrowski et al., 2015; Graf et al., 2019). In addition to AT and SWT, flow discharge has been used 21 to attempt to constrain SWT (Foreman et al., 2001; Tao et al., 2008; St-Hilaire et al., 2011; Grbić et 22 al., 2013; Piotrowski et al., 2015; Graf et al., 2019; Qiu et al., 2020). Traditionally-used model inputs 23 include precipitation (Cole et al., 2014; Jeong et al., 2016; Rozos, 2023), wind direction/speed (Hong 24 and Bhamidimarri, 2012; Cole et al., 2014; Jeong et al., 2016; Kwak et al., 2016; Temizyurek and 25 Dadaser-Celik, 2018; Abdi et al., 2021; Jiang et al., 2022), barometric pressure (Cole et al., 2014), 26 landform attributes (Risley et al., 2003; DeWeber and Wagner, 2014; Topp et al., 2023; Souaissi et 27 al., 2023), and many more (see Table S1).

28 In the last few years, including the day-of-year as an input, DOY (Qiu et al., 2020; Heddam et al., 29 2022: Drainas et al., 2023: Rahmani et al., 2023) and humidity (Cole et al., 2014: Hong and 30 Bhamidimarri, 2012; Kwak et al., 2016; Temizyurek and Dadaser-Celik, 2018; Abdi et al., 2021), 31 have also shown to better capture the seasonal patterns of SWT (Qiu et al., 2020; Philippus et al., 32 2024). With improved access to remote sensing data, there has also been a notable increase of satellite 33 products such as estimates of sky cover (Cole et al., 2014), solar radiation (Kwak et al., 2016; Topp et al., 2023; Majerska et al., 2024), sunshine per day (Drainas et al., 2023) and potential ET (Rozos, 34 35 2023; Topp et al., 2023). However, more research is needed to better understand the influence of 36 newer model inputs on SWT (Zhu and Piotrowski, 2020).

37 Recently, SWT studies focused on the CONUS-scale have chosen to use as many model inputs as 38 available, with Wade et al. (2023), a point-scale CONUS ML study using over 20 variables, while 39 Rahmani et al. (2023) created a LSTM model and considered over 30 variables to simulate SWT. 40 Despite the use of diverse data, the models in these studies performed only satisfactorily and were 41 deemed not generalizable, leaving much room for improvement in CONUS-scale modeling of SWT. 42 With the compilation of larger and larger datasets, feature importance in ML, that is the process of 43 using techniques to assign a score to model input features based on how good the features are at 44 predicting a target variable, can be an efficient way to improve data comprehension, model 45 performance, and model interpretability, the latter of which can dually serve as a transparency marker 46 of which features are driving predictions. Methods for measuring feature importance include using 47 correlation criteria (Pearson's r, Spearman's rho), permutation feature importance (shuffling feature 48 values, measuring decrease in model performance), linear regression feature importance (larger 49 absolute values indicate greater importance), or if using CART/RF/gradient boosting, entropy 50 impurity measurements can be insightful (Venkateswarlu and Anmala, 2023).

51 52

Moved part of section 2.3.1, original (lines 246-253) to section 2.4.2 Model Inputs for ML-SWT (moved to lines 517-523):

3

4 For example, one technique that can be used to improve ML model parameter selection is the 5 Least Absolute Shrinkage and Selection Operator (LASSO), a regression technique used for feature 6 selection (Tibshirani, 1996). Research utilizing ML models for SWT frequency analysis at ungaged 7 basins used the LASSO method to select explanatory variables for two ML models (Souaissi et al., 8 2023). The LASSO method consists of a shrinkage process where the method penalizes coefficients 9 of regression variables by minimizing them to zero (Tibshirani, 1996). The number of coefficients set 10 to zero depends on the adjustment parameter, which controls the severity of the penalty. Thus, the method can perform both feature selection and parameter estimation, an advantage when examining 11 12 large datasets (Xu & Liang, 2021).

13 14

4. In the second section, the authors do an excellent job summarizing model evaluation metrics. However,
considering that ML models are often optimized to achieve superior performance on these metrics, there
is (always) a risk of overfitting. Thus, beyond focusing on metrics, the review should also highlight the
importance of more rigorous evaluation to further assess generalization ability. For instance, if a SWT
model is built to run climate change scenarios, additional testing and more rigorous designs are essential
to evaluate the model's ability to generalize over time. For robust long-term predictions, the model is
supposed to maintain robust predictive performance in completely unseen periods, rather than being

22 limited to a specific temporal range.

AUTHOR RESPONSE: We agree. This comment has similar themes to our response to #3a regarding
 overfitting and highlighting the need for generalization, please see comment #3a for a full response.

For the comment regarding having ML undergo more rigorous testing, we propose adding the following
discussion for more rigorous testing for MLs. We added a few sentences (blue is new) to the Discussion
subsection 4.3 "ML as Knowledge Discovery" where we urge for TUURTs (Temporal, Unseen, Ungaged
Region Tests)' (*revised lines 914-925*):

31 While it is understandable that not every ML-SWT paper aims to explain physical processes, the 32 SWT community should agree on a baseline of tests that all ML-SWT models undergo to assess model 33 robustness and transferability. Specifically, we urge use of TUURTs (Temporal, Unseen, Ungaged 34 Region Tests) for future ML-SWT models as a helpful step towards better modeling practices, 35 increased model transparency and robustness (Fig.1). As stated in figure 1, for TUURTs, testing for 36 "unseen" cases means testing only within the developmental dataset, whereas testing for "ungaged" 37 cases means testing for new sites that have no data and have not been previously seen by the model at 38 all. Due to the difficulty of conducting spatial tests compared to temporal tests, few ML-SWT studies 39 have applied one or two of the tests, and rarely all three (Topp et al., 2023; Hani et al., 2023, Souassi et 40 al., 2023). For example, Siegel et al. (2023), a non-ML SWT paper, tested for ungaged regions and 41 unseen data but did not perform a temporal test. To our knowledge, Philippus et al. (2024), appears to 42 be the only published SWT-ML study that applied TUURTs with some success. We further encourage 43 modelers to shift towards more generalizable models, which are in theory, more capable of performing 44 well across diverse scenarios and datasets and stand to become increasingly important with the 45 unpredictability of climate extremes.

46

30

47 Overall, this review is informative and well-researched, and with more refined organization and deeper
48 exploration of these key issues, it could make a substantial contribution to the field of SWT research.

49 AUTHOR RESPONSE: Thank you! This would certainly not be possible without the insightful

50 feedback from referees.

1 <u>Referee #3 Comments</u>

2 I believe that this manuscript is a very useful and extensive methods literature review regarding stream

3 temperature modeling. I would recommend approval with minor revisions to provide additional details

- 4 from the reviewed literature and correct minor writing aspects; I had no problem with the general
- 5 structure/flow or quality.

AUTHOR RESPONSE: We thank the referee for their time and feedback, we believe the manuscript is
better as a result. We address specific referee comments below. Proposed new/edited text is in BLUE.
Revised lines in the track-changes manuscript are indicated by the statement: (*revised lines XXX-XXX*).

10

Section 2.3.3 ("Newer/recent ML algorithms") introduces RNNs, CNNs, and GNNs sufficiently, but it should probably give some description and reference to attention-based transformers. I am not aware of their application to SWT, but they are responsible for broader interest in ML (e.g., ChatGPT, which was cited earlier) and have had mixed success in hydrologic modeling. This class of models seems easily placed as a future direction.

- AUTHOR RESPONSE: We agree. A literature search on Google Scholar at the end of 2024 found no
 publications specifically using attention-based transformers for SWT modeling, but we can add text about
 their potential to section 2.3.3 (*revised lines 456-462*):
- 19

20 Attention-based transformers are a more novel type of deep learning that has led to advancements 21 in natural language processing, in the form of ChatGPT, Microsoft's CoPilot, Google's Gemini and 22 others. Due to their exponential success in the last few years, attention-based transformer models have 23 been used in geological science fields such as oceanography for sea surface temperature prediction 24 (Shi et al., 2024), hydrology for streamflow and runoff prediction (Ghobadi and Kang, 2022; Wei, 25 2023) and remote sensing for streambed land use change classification (Bansal and Tripathi, 2024). As 26 a relatively new AI tool, attention-based transformers have vet to be used for SWT (to our knowledge). 27 but their applications in other geological science fields suggest it is only a matter of time before their 28 use is observed in SWT modeling.

29 30

2. There are some examples of unusual subsection and paragraph formatting. For example, section 1.1 is
one paragraph which is approximately 1 page long. It seems that this is excessively large for one
paragraph and that a named subsection should perhaps be more than just one (regularly sized) paragraph.
Line 201 has another approximately 1-page-long paragraph, this area might be better organized with
another level of subsections rather than fitting the more extensive references of decision trees into 1
paragraph.

AUTHOR RESPONSE: Thank you for pointing this out. We have consulted other published HESS
 articles and it appears that the first paragraph of a section is not indented but the subsequent paragraphs
 are. We have revised the manuscript to follow the Copernicus manuscript template (screenshot below).

1 Section (as Heading 1)

Suspendisse a elit ut leo pharetra cursus sed quis diam (Smith et al., 2014; Miller and Carter, 2015). Nullam dapibus, ante vitae congue egestas, sem ex semper orci, vel sodales sapien nibh sed lectus. Etiam vehicula lectus quis orci ultricies dapibus. In sit amet lorem egestas, pretium sem sed, tempus lorem.

1.1 Subsection (as Heading 2)

Quisque cursus massa sed urna congue, ac convallis neque consectetur. Proin faucibus neque non metus mollis, suscipit pretium nisl blandit. In hac habitasse platea dictumst. 1

- 2 At the referee's suggestion, we can add subsections to section 2.3.1 to distinguish algorithms as follows: 3 2.3.1.1 K-nearest neighbors (starts line 165), 4 2.3.1.2 Cluster analysis and variants (line 182), 5 2.3.1.3 Support vector machine and regression (line 204), 6 2.3.1.4 Gaussian Process Regression (line 234), 7 2.3.1.5 Decision trees and Classification and Regression Trees (line 255), 8 2.3.1.6 Random Forests and XGBoost (line 272) 9 10 11 We propose making section 2.3.1.6 (original lines 226-253) more concise given the newly separate Model 12 Inputs for ML-SWT section (new section 2.4.2). Below is the revised text for section 2.3.1.6 (blue is new, 13 revised lines 284-313): 14 15 RF and XGBoost, have been used to predictfor daily SWT prediction in 10 for Austrian catchments, 16 Results with results showing ed minor differences in model performance, with a median RMSE 17 difference of 0.08 °C between tested ML models (Feigl et al., 2021). Using RF and XGBoost along 18 with four other ML models, Jiang et al. (2022) tested the performance of six ML models in 19 estimateding daily SWT below dams in China, finding. They found that day of year, was most 20 influential for the prediction of SWT, followed by stream flow flux and AT to be most influential in 21 the prediction of SWT (Jiang et al., 2022). Weierbach et al. (2022) used XGBoost and SVR to predict 22 SWT at monthly time scales for the Pacific Northwest region of the U.S., finding showing that an 23 ensemble XGBoost outperformed all modeling configurations for spatiotemporal predictions in unmonitored basins, In contrast to Jiang et al. (2022), Weierbach et al. (2022) found with AT 24 25 identified as the primary driver of monthly SWT. for all 78 sites in the Pacific Northwest region of 26 the U.S. (which included areas affected by dams), followed by month of year and solar radiation.
- 27 Zanoni et al. (2022) used RF and a deep learning model to develop regional models of SWT and other 28 water quality parameters, finding that with RF performance was comparitively less effective at 29 detecting non-linear relationships than to the deep learning model, though both models identified 30 They found AT to be as most influential, with day of the year, and year of observation as possible 31 replacements where AT was not available (Zanoni et al., 2022).
- 32 Souassi et al. (2023) tested the performance of two ML models, RF and XGBoost, with non-33 parametric models for the regional estimation of maximum SWT at ungaged locations in Switzerland. 34 finding no significant differences between the ML performance and the non-parametric model 35 performances, which was attributed to the lack of a large dataset as required by the ML models. Hani 36 et al. (2023) used four supervised ML models – MARS, GAM, SVM, and RF to model potential 37 thermal refuge area (PTRA) at an hourly timestep for two tributary confluences of the Sainte-38 Marguerite River in Canada. RF had the highest accuracy at both locations in terms of hourly PTRA 39 estimates and modeling SWT (Hani et al., 2023). Wade et al. (2023) conducted a CONUS-scale study 40 using 410 USGS sites with four years of daily SWT and discharge to examine maximum SWT. They 41 used RF to estimate max SWT and thermal sensitivity (Wade et al., 2023), finding that AT was the 42 most influential control followed by other properties (watershed characteristics, hydrology, 43 anthropogenic impact).
- 44

45

46 3. There is an extensive background of traditional ANNs (2.3.2) which is debatably too extensive given the description of ANN variants and backpropagation alternatives (e.g., lines 284-320), which are 47 48 relatively niche and rare. The content already exists and is not wrong, but if length were a concern, I 49 would reduce this area.

50

51 **AUTHOR RESPONSE:** We appreciate the reviewer's feedback and propose making changes to improve the manuscript for readability. Referee #1 made a similar comment about this section, and we propose 52

moving the description of ANN variants and alternatives (lines 263-320) to Appendix A. We think it would still be helpful to keep the ANN information, while also agreeing that it may be too extensive for the main text. In this way, the manuscript can be made more concise while also keeping the details as a section of the manuscript for anyone who is interested in reading further. Following this line of thinking, we removed lines 322-379 from the main text and moved to Appendix A. We provide the following to note the appendix (*revised lines 326-327*):

"For more detail on traditional ANNs, with descriptions of ANN variants and backpropagation alternatives, we refer the reader to appendix A."

We have added Appendix A after the conclusion (*revised lines 1015 – 1075*).

12 13

26

8

9

10 11

4. This work does not address predictive uncertainty, or the lack thereof associated with the ML literature review. I think that would be a worthwhile addition because I suspect most efforts lack that (e.g., referring to <u>https://doi.org/10.5194/hess-26-1673-2022</u>). A counterexample to the lack of uncertainty quantification, which may also be relevant to section 2.5, could be work led by Jacob Zwart focusing on SWT for reservoir operations (thermal releases). Examples being <u>https://doi.org/10.1111/1752-1688.13093</u> or <u>https://doi.org/10.3389/frwa.2023.1184992</u>

AUTHOR RESPONSE: We appreciate the referee's insight in bringing these publications to our
 attention. Based on their relevancy, we have added Klotz et al. 2022, Zwart et al. 2023a and 2023b and
 included their RMSE values in our review. First, we added text on predictive uncertainty in the new
 'Discussion' subsection titled, '4.4 *Future Directions of SWT Modeling*', which also addresses ref #1,
 comment #3, (revised lines 944-991):

27 The utility of ML in hydrologic modeling has advanced significantly, with interest seemingly 28 growing exponentially (Nearing et al., 2021). With the novelty of ML, it is easy to over-value model 29 performance and ignore the physics of the system, but with several decades of ML-experience, we 30 advocate it is necessary to purposefully use ML to address physically-meaningful questions and not 31 just create ML for the sake of creating. Given this, Varadharajan et al. (2022) laid out an excellent 32 discussion on opportunities for advancement of ML in water quality modeling, see section 3 of 33 publication of Varadharajan et al. (2022). (Varadharajan et al., 2022)Here we highlight some of the 34 questions from Varadharajan et al. (2022) that can be considered in the context of what objectives the 35 SWT community should be using in the ML era, namely: 1) How do we use physical knowledge (re: 36 heat exchange equations, radiation influence) to improve models and process understanding? 37 Rahmani et al. (2023) coupled NNs with the physical knowledge from SNTEMP, a one-dimensional 38 stream temperature model that calculates the transfer of energy to or from a stream segment by either 39 heat flux equations or advection, but found that even with SNTEMP, their flexible NNs exhibited 40 substantial variance in prediction and needed to be constrained by further multi-dimensional 41 assessments (Rahmani et al., 2023). In short, if our use of physics in machine learning makes our 42 models worse, we should understand why.

43 A second question that needs addressing is 2) How do we deal with predictive uncertainty in ML 44 used for SWT modeling? According to Moriasi et al. (2007), uncertainty analysis is the process of 45 quantifying the level of confidence in any given model output based on five guidelines: 1) the quality 46 and amount of observations (data), 2) the lack of observations due to poor or limited field monitoring, 47 3) the lack of knowledge of physical processes or operational procedures (instrumentation), 4) the 48 approximation of our mathematical equations, and 5) the robustness of model sensitivity analysis and 49 calibration. For example, in rainfall-runoff modeling, researchers have proposed benchmarking to examine uncertainty predictions of ML rainfall-runoff modeling (Klotz et al., 2022). For stream 50 51 temperature modeling, researchers have attempted to address the role of uncertainty in deep learning 52 model (RGCN, LSTM) predictions using the Monte Carlo Dropout (Zwart, Oliver, et al., 2023) and a unimodal mixture density network approach (Zwart, Diaz, et al., 2023). 53

1 Other questions that SWT-ML studies should consider is 3) How do we make ML models 2 generalize better, specifically with regards to ungaged basins? And 4) How can ML models be 3 improved to predict extremes? As ML models advance to use satellite data, include more sensor 4 networks and/or couple with climate models, there is a logical next step toward creating generalizable 5 models that can account for extremes. In our review, only two papers by the same group (Rahmani et 6 al., 2020, 2023) conducted a CONUS-scale approach towards SWT-ML modeling, omitting 7 hydrologically important regions in the southwest (CA) and southeast (FL). Recently, a satellite 8 remote sensing paper used RF to model monthly stream temperature across the CONUS and tested for 9 temporal (walk-forward validation), unseen and 'true' ungaged regions (Philippus et al., 2024). We 10 have also learned that ML models such as LSTMs, generally only make predictions within the bounds 11 of their training data (Kratzert et al., 2019), which is a limitation for predicting extremes. Thus, we 12 strongly urge the community to work towards ML models that generalize better and/or are more 13 robust towards predictions of extremes.

14 Finally, 5) How can we build ML models such that they are seen as trustworthy and interpretable 15 by the hydrologic community? To answer this question, we must address a technical barrier (black-16 box issues, data limitations, model uncertainty) and a social barrier (i.e., educated skepticism of ML 17 due to novelty, little understanding of computer science basics and/or coding experience). If we are to 18 incorporate ML into decision-making processes, it makes sense that ML must be transparent and 19 understandable to more than just computer or data scientists (Varadharajan et al., 2022). For example, 20 Topp et al. (2023) recently used explainable AI to elucidate how ML architectures affected the SWT 21 model's spatial and temporal dependencies, and how that in turn affected the model's accuracy. 22 Addressing this technical barrier can also be done by improving access to data, which has seen 23 remarkable progress thanks to web repositories such as NSF-funded CUAHSI's Hydro share 24 (CUAHSI, 2024) and GitHub (GitHub, 2024). In the United States, data access to state and locally-25 based data remains limited, and should be addressed. In terms of the social barrier, education about 26 ML and ML-use is key. Societal interest in ML has thankfully also lead to a plethora of educational 27 resources and ML walk-through videos and tutorials in Tensorflow (Abadi et al., 2016), PvTorch 28 (Paszke et al., 2019), and Google Colab (Bisong, 2019). With the speed at which ML-use is evolving, 29 short communication pieces (Lapuschkin et al., 2019) and opinion pieces (Kratzert et al., 2024) with 30 clear examples about an ML-issue and practical solutions will also help make ML challenges more 31 transparent and therefore accessible to the hydrologic community-at-large.

32

33 Added citations used for new subsection, 4.4 Future Directions of SWT Modeling:

- 34 1) Apaydin, H., Taghi Sattari, M., Falsafian, K., and Prasad, R.: Artificial intelligence modelling integrated with
 35 Singular Spectral analysis and Seasonal-Trend decomposition using Loess approaches for streamflow
 36 predictions, Journal of Hydrology, 600, 126506, https://doi.org/10.1016/j.jhydrol.2021.126506, 2021.
- 37 2) Baydaroğlu, Ö. and Demir, I.: Temporal and spatial satellite data augmentation for deep learning-based rainfall nowcasting, Journal of Hydroinformatics, 26, 589–607, https://doi.org/10.2166/hydro.2024.235, 2024.
- 3) CUAHSI. 2024. Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI) Water
 40 Data Portal: https://www.cuahsi.org/community/water-data-portals, last access: 13 November 2024.
- 41 4) Kratzert, F., Gauch, M., Klotz, D. and Nearing, G., 2024. HESS Opinions: Never train an LSTM on a single basin. Hydrology and Earth System Sciences Discussions, 2024, pp.1-19.
- 5) Kwak, J., St-Hilaire, A., and Chebana, F.: A comparative study for water temperature modelling in a small basin, the Fourchue River, Quebec, Canada, Hydrological Sciences Journal, 1–12, https://doi.org/10.1080/02626667.2016.1174334, 2016.
- 6) Philippus, D., Sytsma, A., Rust, A., and Hogue, T. S.: A machine learning model for estimating the temperature of small rivers using satellite-based spatial data, Remote Sensing of Environment, 311, 114271, https://doi.org/10.1016/j.rse.2024.114271, 2024.
- 7) Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C., and Gupta, H. V.:
 What Role Does Hydrological Science Play in the Age of Machine Learning?, Water Resources Research, 57, e2020WR028091, https://doi.org/10.1029/2020WR028091, 2021.

- 8) Skoulikaris, C., Venetsanou, P., Lazoglou, G., Anagnostopoulou, C., and Voudouris, K.: Spatio-Temporal Interpolation and Bias Correction Ordering Analysis for Hydrological Simulations: An Assessment on a Mountainous River Basin, Water, 14, 660, https://doi.org/10.3390/w14040660, 2022.
- 9) Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: A Simple Way to
 Prevent Neural Networks from Overfitting, Journal of Machine Learning Research, 15, 30, 2014.
- 6 10) Yang, M., Yang, Q., Shao, J., Wang, G., and Zhang, W.: A new few-shot learning model for runoff prediction:
 7 Demonstration in two data scarce regions, Environmental Modelling & Software, 162, 105659, https://doi.org/10.1016/j.envsoft.2023.105659, 2023.
- 9 11) GitHub. 2024. About Git and Github: https://docs.github.com/en/get-started/start-your-journey/about-github 10 and-git, last access: 14 November 2024.
- 12) Lapuschkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W. and Müller, K.R., 2019. Unmasking Clever
 Hans predictors and assessing what machines really learn. Nature communications, 10(1), p.1096.
- 13) Zwart, J.A., Oliver, S.K., Watkins, W.D., Sadler, J.M., Appling, A.P., Corson-Dosch, H.R., Jia, X., Kumar, V. and Read, J.S., 2023. Near-term forecasts of stream temperature using deep learning and data assimilation in support of management decisions. JAWRA Journal of the American Water Resources Association, 59(2), pp.317-337.
- 16 14) Zwart, J.A., Diaz, J., Hamshaw, S., Oliver, S., Ross, J.C., Sleckman, M., Appling, A.P., Corson-Dosch, H., Jia,
 X., Read, J. and Sadler, J., 2023. Evaluating deep learning architecture and data assimilation for improving
 water temperature forecasts at unmonitored locations. *Frontiers in Water*, 5, p.1184992.
- 15) Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S. and
 Nearing, G., 2022. Uncertainty estimation with deep learning for rainfall-runoff modeling. Hydrology and
 Earth System Sciences, 26(6), pp.1673-1693.
- 16) M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, R. Jozefowicz, Y. Jia, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, M. Schuster, R. Monga, S. Moore, D. Murray, C. Olah, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems. 2015. TensorFlow. Website: https://www.tensorflow.org/
- 17) A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. Website: https://arxiv.org/abs/1912.01703
- 32 18) Bisong, E. (2019). Google Colaboratory. In: Building Machine Learning and Deep Learning Models on Google
 33 Cloud Platform. Apress, Berkeley, CA. Website: <u>https://doi.org/10.1007/978-1-4842-4470-8_7</u>
- 34

- 38 Further focusing on the Delaware River Basin, Zwart, Oliver, et al. (2023) used data assimilation 39 and an LSTM to generate 1-day and 7-day forecasts of daily maximum SWT for the purpose of aiding 40 reservoir managers in decisions about when to release water to cool streams. Following up on this 41 study was Zwart, Diaz, et al. (2023), who used a LSTM and a RGCN, to generate 7-day forecasts of 42 daily maximum SWT for monitored and unmonitored locations in the Delaware River Basin. The 43 study found that the RGCN with data assimilation performed best for ungaged locations and for 44 higher SWT, which can serve as valuable information for reservoir operators to consider while 45 drafting release schedules.
- 46 47
- 48 5. In section 3 (e.g., 3.1, 3.3, 3.4), I would recommend adding some discussion regarding the equivalence
 49 or lack of between lower-case r and r-squared, upper-case R-squared, and NSE. I am very comfortable

^{We have also added text to section 2.5 Decision Support with the provided citations (}*revised lines 663-668*):
37

2 are equivalent, but I am less comfortable making the assertation that lower case r and r-squared are (in all 3 the papers reporting this value). This is likely further complicated by the reviewed literature using the 4 lower-case r-squared and R-squared interchangeably, but given the 0-1 range, the high value skew, and 5 the special case/conditional equivalences, I believe these values should all be reported together to 6 characterize goodness of fit – especially that upper case R-squared and NSE should not be separated. 7 8 **AUTHOR RESPONSE**: We agree. We propose the following to address the referee's comments: 9 Revise section 3.1 text to distinguish between lower-case r, r-squared r^2 , and upper-case R^2 (revised 10 11 lines 697-710): 12 13 Pearson's r, also known as the correlation coefficient, is useful for determining the strength and 14 direction (i.e., positive, negative) of a simple linear relationship (Helsel and Hirsch, 2002). Values of r, 15 range from -1 to +1, where r < 0 indicates a negative correlation and r > 0 indicates a positive 16 correlation (Legates and McCabe, 1999). The square of r is denoted as $\mathbb{R}^2 r^2$, or known as the square of 17 the correlation coefficient, with values of r^2 ranging from 0 to 1. The r^2 metric is commonly used in 18 simple linear regression to assess the goodness of fit by of determination, which represents measuring 19 the fraction of the variance in one variable (i.e., observations) that can be explained by the other 20 variable (i.e., predictors). The metric r^2 tends to be confused with R^2 , the latter which is a statistical 21 measure that represents the proportion of variance explained by the independent variable(s) in a multiple linear regression model (Helsel and Hirsch, 2002). Part of the confusion may be related to the 22 fact that R^2 shares the same range of from 0 to 1, with $R^2 = 1$ suggesting indicating that the model can 23 24 explain all the variance, and vice versa. We note that while both r^2 and R^2 share similarities in that they 25 measure the proportion of variance, R^2 is more commonly used for multiple linear regression context, 26 while r^2 is best suited for simple linear regressions. To reduce confusion, we strongly suggest that r, r^2 27 and R^2 always be reported together (even if as a supplement to a manuscript) to characterize goodness-28 of-fit. The r and R^2 metrics are typically used for normally distributed data that follows a bivariate-29 normal distribution (Helsel and Hirsch, 2002). 30 31 Add text stating that upper R^2 and NSE should always be provided together in section 3.4: 32 -33 34 1st paragraph, added after 1st sentence (*revised lines* 785-789): 35 Having reviewed the literature and in agreement with previous published recommendations 36 (Moriasi et al., 2007), we recommend that a combination of standard regression (i.e., r, r^2, R^2). 37 dimensionless (i.e., NSE), and error index statistics (i.e., RMSE, MAE, PBIAS) be used for model 38 evaluation and reported together in future publications. 39 40 2nd paragraph, remove the statement about r and R2 (*revised lines 792-793*): 41 We note that for the 11 studies that used Pearson's r (see Table S1), and given that r and R2 are 42 directly related, we converted r to R2 for ease of comparison on fig. 1. 43 44 3rd paragraph, added last sentence (*revised lines 803-805*): 45 Overall, these complimentary metrics should always be reported together as they provide a 46 broader evaluation of model performance, i.e., NSE measures a model's predictive skill and error 47 variance, while R^2 assesses how well the model explains the variability of the data. 48 In section 3.4, remove all r^2 values from Figure 1, only R^2 citations (17) remain. The median R^2 for 49 training stayed the same (0.93), while the testing R^2 went from 0.95 to 0.94, and the validation R^2 went 50 51 from 0.92 to 0.93. Overall, changes were insignificant. Below is a screenshot of the "Original (top)" 52 and "Revised (bottom)" Figure 1 for reference.

stating that for the purpose of this continuously valued model evaluation, upper case R-squared and NSE



- 37

38 6. In line 761, it feels controversial and a step too far to say ML models should be held to a higher 39 standard. It feels less problematic to apply these higher, seemingly attainable standards to all SWT models. For example, a physics-based model is not "very good" by virtue of being a physics-based model, 40 41 instead it is the same "satisfactory" label because its physics are not sufficient or accurate enough to do

42 what the ML models can.

43 **AUTHOR RESPONSE**: We appreciate the referee's point of view. Perhaps instead of saying "separate, 44 higher standard", we can say "additional standards" (see revised line 1002), but we think that additional 45 standards are warranted nonetheless, not only in terms of performance metrics but also to improve model 46 transparency, eradicate black-box confusion and encourage user confidence. We disagree that a physics-47 based model should be in the same "satisfactory" performance metric category because the intention of 48 performance metrics is to identify what fits the data best (which data-driven ML excel at), whereas the 49 general intention of physics-based models is to adhere to whatever governing equations have been 50 employed. Our review indicates that we have been somewhat blinded by the excellence of ML 51 performance metrics relative to physics-based and statistically-based models, and more awareness is 52 needed moving forward.

1

7. If possible, in addition to considering spatial extents and temporal resolution of the papers, it would be interacting to know the aggregation level of data if that is reported and what all the pageibilities are. For

- interesting to know the aggregation level of data if that is reported and what all the possibilities are. For
 example, individual gages with input data collected at the same gage location in situ, remotely sensed data
- example, individual gages with input data collected at the same gage location in situ, remotely sensed data
 subset to the drainage area for the reach that a gage is on. Are any works modeling dense transects along a
- 6 river or modeling raster grid cells up and across a river (i.e., the 2D surface area), etc.

AUTHOR RESPONSE: Thank you for the opportunity to clarify. We provided supplementary table S1
 to summarize study information regarding time period, temporal resolution, spatial resolution and

9 hydrometeorological parameters considered by the cited studies. Responding to your comment, in our

10 review, we saw that the aggregation level of data is more often than not, left unreported and unclear by

11 studies (and reporting is not mandatory as a lot of data is pre-processed before utilization in modeling,

- 12 adding to transparency questions). We do think discerning all the possibilities of data aggregation could
- make for an interesting follow-up study for the larger hydrologic community, which could focus solely ondata manipulation, processing and augmentation for ML.
- 14 data ma 15
- 16

17 Additional literature to consider. Not necessary

18 8. The paragraph at line 385 related to process guidance prompted me to recommend

19 https://doi.org/10.1029/2023WR035327 as very relevant. The reference is concerned with comparing

20 different hybrid ML methods for SWT modeling to represent groundwater processes which aren't as

21 represented here (e.g., relative to reservoir influence/reservoir adjacent modeling).

22 AUTHOR RESPONSE: Thank you for the suggestion, we agree that the challenge of including

23 groundwater influence in SWT modeling warrants more research. We want to clarify that we did not

include this reference as it appears to be a conference paper and not subjected to journal standards of peer

review. That being said, the authors of the suggested manuscript went on to publish similar work in Water

- Resources Research, which we cite in this review (Topp et al., 2023).
- 27

9. In section 4.2, https://doi.org/10.1029/2020WR028091 may be a very relevant addition in-line with the
 author's narrative.

AUTHOR RESPONSE: Thank you for the suggestion, we enjoyed reading it and think it insightful. We
 added it to new 'Discussion' subsection, titled '4.4 Future Directions of SWT Modeling', in the first
 sentence (please see our response to ref #1, comment #4 for the full text), *revised lines 944-945*:

"The utility of ML in hydrologic modeling has come a long way, with interest seemingly growing exponentially (Nearing et al., 2021)."

35 36 37

34

38 Minor writing comments:

39 1.The sentence beginning on line 51 perhaps uses too bold language when stating "AI ... create

40 reasonable choices". Many users of AI and scientists have concerns regarding the reasonableness of AI.

41 Maybe it would be more accurate to further connect with the latter part of that sentence and say that "AI

42 ... learn optimal patterns to meet stated objectives" (which may or may not be broadly reasonable)

43 AUTHOR RESPONSE: That is a good point. Reasonableness is fluid. We agree with the referee and
44 have updated the sentence as follows (*revised lines 51-53*):
45

"Artificial intelligence (AI) describes technologies that can incorporate and assess inputs from an
 environment, create reasonable choices, learn optimal patterns and implement actions to meet stated

1 2 3 4 5	objectives or performance metrics (Xu & Liang, 2021; Varadharajan et al., 2022)." 2. Starting at line 131, "We define newer ML as those introduced in hydrologic modeling in the few years," perhaps this should say "in recent years"?
6 7 8 9	AUTHOR RESPONSE: We agree, thank you for the suggestion, we have updated the text to say, "in recent years" (<i>revised line 159</i>).
10	3. At line 380, although it can be inferred, "WNN" is never explicitly defined.
11 12 13	AUTHOR RESPONSE: Thank you for catching that, we have defined the acronym (<i>revised line 445</i>).
14	4. At line 541, "all journals examined used least one", perhaps this should say, "at least one"
15 16	AUTHOR RESPONSE: Thank you! We have added the word "at" (<i>revised line 681</i>).
17 18 19	5. By typo/mistake, it appears that two subsections in section 3 are titled "Model Performance Metrics: Error Indices"
20 21 22 23 24	AUTHOR RESPONSE: Yes, thank you for catching this error. Subsection 3.3 should have said "Dimensionless" because the subsection summarizes dimensionless metrics. We have updated the subsection header accordingly (<i>title, revised line 739</i>).
25	6. At line 610, there is a typo claiming an upper bound of -1
26 27 28 29	AUTHOR RESPONSE: Yes, that was a typo. Thank you for catching that, we have updated the text to just say "0 to 1" (<i>revised line 759</i>).
30 31 32 33 34	7. I have the benefit of reviewing 3rd, so I read the other reviewer's comments after making my own. I agree that a characterization of the validation and test sets used would be very beneficial (e.g., spatial, temporal, spatiotemporal exclusion, etc.), but I believe the concerns of overfitting are potentially overstated by the other reviewers given that this manuscript reports train, validation, and test set metrics (and the very strong agreement between the three).
35 36 37 38 39	AUTHOR RESPONSE: Thank you for your time and energy in reviewing this manuscript. With regards to the concerns of overfitting, we include below our response to referee #1, comment #1A. We think that the referee comment with regard to "characterization of the validation and test sets" is related to referee #1, comment #1B, which we also include below:
40 41	New subsection 2.4.1 Identifying Model Complexity (revised lines 464-483):
4 1 42 43 44 45 46 47	The strong success of ML-use in SWT modeling warrants a brief and broad overview on identifying model complexity to minimize overfitting and underfitting" of models. When a model is too complex, i.e., has too many features or parameters relative to the number of observations, or is forced to overextend its capabilities, i.e., make predictions with insufficient training data, the model runs the risk of overfitting (Srivastava et al., 2014). An overfitted model fits the training data "too well", capturing noise and details that provide high accuracy on a training dataset, only to perform poorly
48	once the model encounters "unseen" data in testing/validation (Xu and Liang, 2021). Scenarios where

overfitting may be temporarily acceptable are: 1) model development is at preliminary stages, the
interest is in a "proof of life" concept, 2) when the objective is to identify heavily-relied on features
by the model, i.e., feature importance, or 3) in highly-controlled modeling environments where the
expected data will be consistently similar to the training dataset. The latter is more likely in industrial
applications and unlikely in the changing nature of hydrology.

7 In contrast, underfitting occurs when a model is too simple to capture any patterns in the data, 8 which can also lead to unsatisfactory performance in training, testing and validation. Underfitting can 9 occur with inadequate model features, poor model complexity or when regularization techniques, 10 (e.g., L1 or L2 regularization), are over-used, making the model too rigid and unable to respond to changes in the data. Given the propensity of ML models to effectively learn the training data, 11 12 underfitting is less of an issue in ML whereas overfitting can be widespread. In Figure 1, we present 13 an example workflow that researchers can use to transition away from overfitting and towards 14 generalizability. In the five-step outline (Fig. 1), we suggest the need for "Temporal, Unseen, 15 Ungaged Region Tests" (TUURTs), which is a call for temporal and spatially-focused testing that can 16 be used to strengthen model robustness.

18 Response to ref #1, comment #1B: We have added a few sentences (blue is new) to the Discussion
19 subsection 4.3, "ML Use for Knowledge Discovery" where we further urge for the use of TUURTs
20 (Temporal, Unseen, Ungaged Region Tests)' (*revised lines 914-925*):

22 While it is understandable that not every ML-SWT paper aims to explain physical processes, the 23 SWT community should agree on a baseline of tests that all ML-SWT models undergo to assess 24 model robustness and transferability. Specifically, we urge use of TUURTs (Temporal, Unseen, 25 Ungaged Region Tests) for future ML-SWT models as a helpful step towards better modeling 26 practices, increased model transparency and robustness (Fig.1). As stated in figure 1, for TUURTs, 27 testing for "unseen" cases means testing only within the developmental dataset, whereas testing for 28 "ungaged" cases means testing for new sites that have no data and have not been previously seen by 29 the model at all. Due to the difficulty of conducting spatial tests compared to temporal tests, few ML-30 SWT studies have applied one or two of the tests, and rarely all three (Topp et al., 2023; Hani et al., 31 2023, Souassi et al., 2023). For example, Siegel et al. (2023), a non-ML SWT paper, tested for 32 ungaged regions and unseen data but did not perform a temporal test. To our knowledge, Philippus et 33 al. (2024), appears to be the only published SWT-ML study that applied TUURTs with some success. 34 We further encourage modelers to shift towards more generalizable models, which are in theory, 35 more capable of performing well across diverse scenarios and datasets, and stand to become 36 increasingly important with the unpredictability of climate extremes.

37

17

38 Disclaimer from Reviewer: I propose some additional literature (n = 4-5), and I am a coauthor on 1 of

39 them. I do not view including that literature as mandatory, and only proposed additional sources based

40 *on their relevance to the content of this manuscript. I selected "No" to anonymity to avoid any*

41 *appearance of subversive influence.*

42 AUTHOR RESPONSE: Thank you! This would certainly not be possible without the insightful
 43 feedback from referees.