

1 **Referee #3 Comments**

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.

6 **AUTHOR RESPONSE:** We thank the referee for their time and feedback, we believe the manuscript is
7 better as a result. We address specific referee comments below. **Proposed new/edited text is in BLUE.**

8
9

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

15 **AUTHOR RESPONSE:** We agree. A literature search on Google Scholar in November 2024 found no
16 publications specifically using attention-based transformers for SWT, but we are happy to add some text
17 about their potential to section 2.3.3:

18

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

28

29

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

36 **AUTHOR RESPONSE:** We appreciate the opportunity to clarify. For section 1.1 (line 35), the 2nd
37 paragraph begins on line 46, with the words “*Aided by the continued...*”. The same occurs after Line 201,
38 where the RF and XGBoost paragraph begins on line 238. The manuscript follows the Copernicus
39 manuscript template (screenshot below) which appears to not provide for paragraph indentation.

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 138),
- 4 2.3.1.2 Cluster analysis and variants (line 145),
- 5 2.3.1.3 Support vector machine and regression (line 160),
- 6 2.3.1.4 Gaussian Process Regression (line 189),
- 7 2.3.1.5 Decision trees and Classification and Regression Trees (line 202),
- 8 2.3.1.6 Random Forests and XGBoost (line 215)

9
10 We also think that we can make section 2.3.1.6 (lines 226-253) more concise now because model inputs
11 are now a separate section (section 2.4.X). Below is our suggested reduction, with the last LASSO
12 paragraph also being moved to model inputs and selection:

13
14 ~~Feigl et al. (2021) tested the performance of six ML models, including RF and XGBoost have been~~
15 ~~used to predict SWT for Austrian catchments with minor differences in model performance, for daily~~
16 ~~SWT prediction in 10 Austrian catchments. Results showed minor difference in model performance,~~
17 with a median RMSE difference of 0.08 °C between tested ML models (Feigl et al., 2021). Using RF
18 and XGBoost along with four other ML models, Jiang et al. (2022) ~~tested the performance of six ML~~
19 ~~models in estimating estimated~~ daily SWT below dams in China, finding. ~~They found that day of~~
20 ~~year, stream flow flux and AT to be was~~ most influential for the prediction of SWT, ~~followed by~~
21 ~~stream flow flux and AT~~ (Jiang et al., 2022). Weierbach et al. (2022) used XGBoost and SVR to
22 predict SWT at monthly time scales for the Pacific Northwest region of the U.S., finding that an
23 ensemble XGBoost outperformed all modeling configurations for spatiotemporal predictions in
24 unmonitored basins. ~~In contrast to Jiang et al. (2022), Weierbach et al. (2022) found AT as the~~
25 ~~primary driver of monthly SWT, for all 78 sites in the Pacific Northwest region of the U.S. (which~~
26 ~~included areas affected by dams), followed by month of year and solar radiation.~~ Zanoni et al. (2022)
27 used RF and a deep learning model to develop regional models of SWT and other water quality
28 parameters, finding that RF performance was comparatively less effective at detecting non-linear
29 relationships, though both models identified AT as most influential ~~than to the deep learning model.~~
30 ~~They found AT to be 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 compared 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 thermal
37 refuge area (PTRA) at an hourly timestep for two tributary confluences of the Sainte-Marguerite River
38 in Canada. RF had the highest accuracy at both locations in terms of hourly PTRA estimates and
39 modeling SWT (Hani et al., 2023). Wade et al. (2023) conducted a CONUS-scale study using RF 410
40 USGS sites with four years of daily SWT and discharge to examine maximum SWT. They used RF to
41 estimate found that max SWT and thermal sensitivity (Wade et al., 2023), finding Study findings
42 identified that AT was the as most influential control followed by other properties (watershed
43 characteristics, hydrology, anthropogenic impact).

44
45
46 **3.** There is an extensive background of traditional ANNs (2.3.2) which is debatably too extensive given
47 the description of ANN variants and backpropagation alternatives (e.g., lines 284-320), which are
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 are open to making changes to
52 improve the manuscript for readability. Referee #1 made a similar comment about this section, and we
53 now propose providing the description of ANN variants and alternatives (lines 263-320) as part of an

1 appendix. We think it would still be helpful to keep the ANN information, but we also agree that it may
 2 be too extensive for the main text. In this way, the manuscript can be made more concise while also
 3 keeping the details as a section of the manuscript for anyone who is interested in reading further.
 4 Following this line of thinking, we can add the following to point the reader to the appendix:

5
 6 “For more detail on traditional ANNs, with descriptions of ANN variants and backpropagation
 7 alternatives, we refer the reader to appendix A.”
 8
 9

10 4. This work does not address predictive uncertainty, or the lack thereof associated with the ML literature
 11 review. I think that would be a worthwhile addition because I suspect most efforts lack that (e.g., referring
 12 to <https://doi.org/10.5194/hess-26-1673-2022>). A counterexample to the lack of uncertainty
 13 quantification, which may also be relevant to section 2.5, could be work led by Jacob Zwart focusing on
 14 SWT for reservoir operations (thermal releases). Examples being [https://doi.org/10.1111/1752-](https://doi.org/10.1111/1752-1688.13093)
 15 [1688.13093](https://doi.org/10.1111/1752-1688.13093) or <https://doi.org/10.3389/frwa.2023.1184992>
 16

17 **AUTHOR RESPONSE:** We appreciate the referee’s insight in bringing these publications to our
 18 attention. Based on their relevancy, we have added Klotz et al. 2022, Zwart et al. 2023a and 2023b to our
 19 manuscript and included their RMSE values in our review. First, we added text on predictive uncertainty
 20 in the new ‘Discussion’ subsection, titled ‘Future Directions of SWT Modeling’ (this section also
 21 addresses ref #1, comment #3):
 22

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

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

50 Other questions that SWT-ML studies should consider is 3) How do we make ML models
 51 generalize better, specifically with regards to ungaged basins? And 4) How can ML models be
 52 improved to predict extremes? As ML models advance to use satellite data, include more sensor
 53 networks and/or couple with climate models, there is a logical next step toward creating generalizable

1 models that can account for extremes. In our review, only two papers by the same group (Rahmani et
 2 al., 2020, 2023) conducted a CONUS-scale approach towards SWT-ML modeling, omitting
 3 hydrologically important regions in the southwest (CA) and southeast (FL). Recently, a satellite
 4 remote sensing paper used RF to model monthly stream temperature across the CONUS and tested for
 5 temporal (walk-forward validation), unseen and ‘true’ ungaged regions (Philippus et al., 2024). We
 6 have also learned that ML models such as LSTMs, generally only make predictions within the bounds
 7 of their training data (Kratzert et al., 2019), which is a limitation for predicting extremes. Thus, we
 8 strongly urge the community to work towards ML models that generalize better and/or are more
 9 robust towards predictions of extremes.

10 Finally, 5) How can we build ML models such that they are seen as trustworthy and
 11 interpretable by the hydrologic community? To answer this question, we must address a technical
 12 barrier (black-box issues, data limitations, model uncertainty) and a social barrier (i.e., educated
 13 skepticism of ML due to novelty, little understanding of computer science basics and/or coding
 14 experience). If we are to incorporate ML into more of the decision-making process, it makes sense
 15 that ML must be transparent and understandable to more than just computer scientists (Varadharajan
 16 et al., 2022). For example, Topp et al. (2023) recently used explainable AI to elucidate how ML
 17 architectures affected the SWT model’s spatial and temporal dependencies, and how that in turn
 18 affected the model’s accuracy. Addressing this technical barrier can also be done by improving access
 19 to data, which has seen remarkable progress thanks to web repositories such as NSF-funded
 20 CUAHSI’s Hydro share (CUAHSI, 2024) and GitHub (GitHub, 2024). In the United States, data
 21 access to state and locally-based data remains limited, and should be addressed. In terms of the social
 22 barrier, education about ML and ML-use is key. Societal interest in ML has thankfully also lead to a
 23 plethora of educational resources and ML walk-through videos and tutorials in Tensorflow (Abadi et
 24 al., 2015), PyTorch (Abadi et al., 2015), and Google Colab (Bison, 2019). With how fast ML-use is
 25 evolving, short communication pieces (Lapuschkin et al., 2019) and opinion pieces (Kratzert et al.,
 26 2024) with clear examples about an ML-issue and practical solutions could also help make ML
 27 challenges more transparent and therefore accessible to the hydrologic community-at-large.

28
 29 We have added a few lines to section 2.5 Decision Support with the provided citations:

30
 31 Further focusing on the Delaware River Basin, Zwart, Oliver, et al. (2023) used data assimilation
 32 and an LSTM to generate 1-day and 7-day forecasts of daily maximum SWT for the purpose of aiding
 33 reservoir managers in decisions about when to release water to cool streams. Following up on this
 34 study was Zwart, Diaz, et al. (2023), who used a LSTM and a RGCN, to generate 7-day forecasts of
 35 daily maximum SWT for monitored and unmonitored locations in the Delaware River Basin. The
 36 study found that the RGCN with data assimilation performed best for ungaged locations and for
 37 higher SWT, which can serve as valuable information for reservoir operators to consider while
 38 drafting release schedules.

39
 40 5. In section 3 (e.g., 3.1, 3.3, 3.4), I would recommend adding some discussion regarding the equivalence
 41 or lack of between lower-case r and r -squared, upper-case R -squared, and NSE. I am very comfortable
 42 stating that for the purpose of this continuously valued model evaluation, upper case R -squared and NSE
 43 are equivalent, but I am less comfortable making the assertion that lower case r and r -squared are (in all
 44 the papers reporting this value). This is likely further complicated by the reviewed literature using the
 45 lower-case r -squared and R -squared interchangeably, but given the 0-1 range, the high value skew, and
 46 the special case/conditional equivalences, I believe these values should all be reported together to
 47 characterize goodness of fit – especially that upper case R -squared and NSE should not be separated.

48
 49 **AUTHOR RESPONSE:** We agree. We propose the following to address the referee’s comments:

- 50
 51 - Revise section 3.1 text to clearly distinguish between lower-case r , r -squared r^2 , and upper-case R^2 :

52
 53 The square of r is denoted as r^2 , or known as the square of the correlation coefficient, with values

of r^2 ranging from 0 to 1. The r^2 metric is commonly used in simple linear regression to assess the goodness of fit by measuring the fraction of the variance in one variable (i.e., observations) that can be explained by the other variable (i.e., predictors). The metric r^2 tends to be confused with R^2 , the latter which is a statistical 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 fact that R^2 shares the same range of 0 to 1, with $R^2 = 1$ indicating that the model can explain all the variance, and vice versa. We note here that while both r^2 and R^2 share similarities in that they measure the proportion of variance, R^2 is more commonly used for multiple linear regression context, while r^2 is best suited for simple linear regressions. To prevent confusion, we strongly suggest that r , r^2 and R^2 always be reported together (even if as a supplement to a manuscript) to characterize goodness-of-fit. The r and R^2 metrics are typically used for normally distributed data that follows a bivariate normal distribution (Helsel and Hirsch, 2002).

- Add text stating that upper R^2 and NSE should always be provided together in section 3.4:

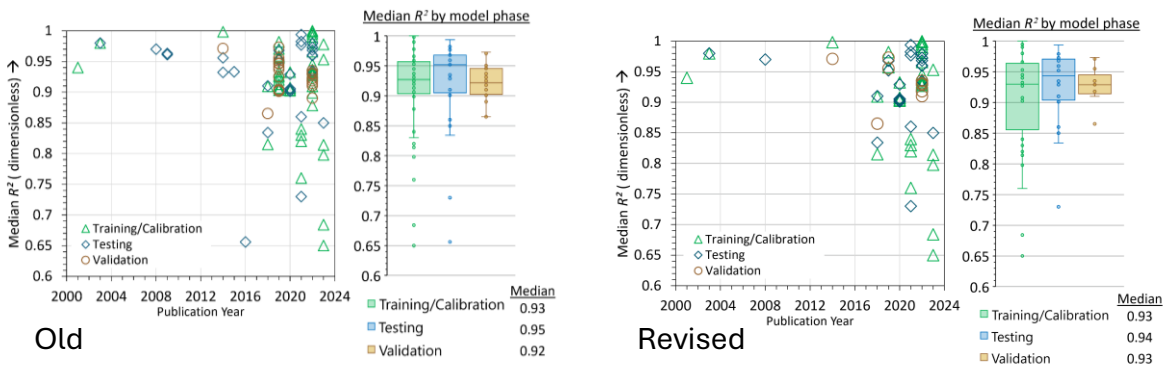
1st paragraph, added after 1st sentence:

Having reviewed the literature and in agreement with previous published recommendations (Moriassi et al., 2007), we recommend that a combination of standard regression (i.e., r , r^2 , R^2), dimensionless (i.e., NSE), and error index statistics (i.e., RMSE, MAE, PBIAS) be used for model evaluation and reported together in future publications.

3rd paragraph, added last sentence:

Overall, these complimentary metrics should always be reported together as they provide a broader evaluation of model performance, i.e., NSE measures a model’s predictive skill and error variance, while R^2 assesses how well the model explains the variability of the data.

- In section 3.4, remove all r^2 values from Figure 1, only R^2 citations (17) remain. The median R^2 for training stayed the same (0.93), while the testing R^2 went from 0.95 to 0.94, and the validation R^2 went from 0.92 to 0.93. Overall, changes were insignificant. Below is a screenshot of the “Old (left)” and “Revised (right)” Figure 1 for reference.



6. In line 761, it feels controversial and a step too far to say ML models should be held to a higher 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, instead it is the same "satisfactory" label because its physics are not sufficient or accurate enough to do what the ML models can.

AUTHOR RESPONSE: We appreciate the referee’s point of view and are open to discussion. Perhaps instead of saying “higher standard”, we can say “additional standards”, but we think that additional

1 standards are warranted nonetheless, not only in terms of performance metrics but also to improve model
 2 transparency, eradicate black-box confusion and encourage user confidence. We disagree that a physics-
 3 based model should be in the same “satisfactory” performance metric category because the intention of
 4 performance metrics is to identify what fits the data best (which data-driven ML excel at), whereas the
 5 general intention of physics-based models is to adhere to whatever governing equations have been
 6 employed. This review shows that we have been blinded by the excellence of ML performance metrics
 7 relative to physics-based and statistically-based models, and we need to be aware of this short sight
 8 moving forward.

9
 10
 11 7. If possible, in addition to considering spatial extents and temporal resolution of the papers, it would be
 12 interesting to know the aggregation level of data - if that is reported and what all the possibilities are. For
 13 example, individual gages with input data collected at the same gage location in situ, remotely sensed data
 14 subset to the drainage area for the reach that a gage is on. Are any works modeling dense transects along a
 15 river or modeling raster grid cells up and across a river (i.e., the 2D surface area), etc.

16 **AUTHOR RESPONSE:** Thank you for the opportunity to clarify. We provided supplementary table S1
 17 to summarize study information regarding time period, temporal resolution, spatial resolution and
 18 hydrometeorological parameters considered by the cited studies. Responding to your comment, in our
 19 review, we saw that the aggregation level of data is more often than not, left unreported and unclear by
 20 studies (and reporting is not mandatory as a lot of data is pre-processed before utilization in modeling,
 21 adding to transparency questions). We do think discerning all the possibilities of data aggregation could
 22 make for an interesting follow-up study for the larger hydrologic community, which could focus solely on
 23 data manipulation, processing and augmentation for ML.

24
 25
 26 **Additional literature to consider. Not necessary**

27 8. The paragraph at line 385 related to process guidance prompted me to recommend
 28 <https://doi.org/10.1029/2023WR035327> as very relevant. The reference is concerned with comparing
 29 different hybrid ML methods for SWT modeling to represent groundwater processes which aren't as
 30 represented here (e.g., relative to reservoir influence/reservoir adjacent modeling).

31 **AUTHOR RESPONSE:** Thank you for the suggestion, we agree that the challenge of including
 32 groundwater influence in SWT modeling warrants more research. We want to clarify that we did not
 33 include this reference as it appears to be a conference paper and not subjected to journal standards of peer
 34 review. That being said, the authors of the suggested manuscript went on to publish similar work in Water
 35 Resources Research, which we cite in this review (Topp et al., 2023).
 36

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

39 **AUTHOR RESPONSE:** Thank you for the suggestion, we enjoyed reading it and think it insightful. We
 40 added it to a proposed new ‘Discussion’ subsection, titled ‘Future Directions of SWT Modeling’, in the
 41 first sentence (please see our response to ref #1, comment #4 for the full text):

42
 43 “The utility of ML in hydrologic modeling has come a long way, with interest seemingly growing
 44 exponentially (Nearing et al., 2021).”
 45
 46
 47
 48
 49

1 **Minor writing comments:**

2 1. The sentence beginning on line 51 perhaps uses too bold language when stating “AI ... create
3 reasonable choices”. Many users of AI and scientists have concerns regarding the reasonableness of AI.
4 Maybe it would be more accurate to further connect with the latter part of that sentence and say that “AI
5 ... learn optimal patterns to meet stated objectives” (which may or may not be broadly reasonable)

6 **AUTHOR RESPONSE:** That is a good point. Reasonableness is fluid. We agree with the referee and
7 have updated the sentence as follows:

8
9 “Artificial intelligence (AI) describes technologies that can incorporate and assess inputs from an
10 environment, ~~create reasonable choices~~, **learn optimal patterns** and implement actions to meet stated
11 objectives or performance metrics (Xu & Liang, 2021; Varadharajan et al., 2022).”
12
13

14 2. Starting at line 131, “We define newer ML as those introduced in hydrologic modeling in the few
15 years,” perhaps this should say “in recent years”?

16 **AUTHOR RESPONSE:** We agree, thank you for the suggestion, we have updated the text to say, “**in**
17 **recent years**”.

18
19

20 3. At line 380, although it can be inferred, “WNN” is never explicitly defined.

21 **AUTHOR RESPONSE:** Thank you for catching that, we have defined the acronym.
22
23

24 4. At line 541, “all journals examined used least one”, perhaps this should say, “at least one”

25 **AUTHOR RESPONSE:** Thank you! We have added the word “**at**”.
26
27

28 5. By typo/mistake, it appears that two subsections in section 3 are titled "Model Performance Metrics:
29 Error Indices"

30 **AUTHOR RESPONSE:** Yes, thank you for catching that mistake. Subsection 3.3 should have said
31 “Model Performance Metrics: Dimensionless” because the subsection summarizes NSE, KGE, etc. We
32 have updated the subsection header accordingly.
33
34

35 6. At line 610, there is a typo claiming an upper bound of -1

36 **AUTHOR RESPONSE:** Yes, that was a typo. Thank you for catching that, we have updated the text to
37 just say “**0 to 1**”.
38
39

40 7. I have the benefit of reviewing 3rd, so I read the other reviewer’s comments after making my own. I
41 agree that a characterization of the validation and test sets used would be very beneficial (e.g., spatial,
42 temporal, spatiotemporal exclusion, etc.), but I believe the concerns of overfitting are potentially
43 overstated by the other reviewers given that this manuscript reports train, validation, and test set
44 metrics (and the very strong agreement between the three).

45 **AUTHOR RESPONSE:** Thank you for your time and energy in reviewing this manuscript. With regards
46 to the concerns of overfitting, we include below our response to referee #1, comment #1A. We think that

1 the referee comment with regard to “characterization of the validation and test sets” is related to referee
 2 comment #1B, which we also include below:

3
 4 Section 2.4.X Overfitting and Underfitting:

5
 6 When a model is too complex, i.e., has too many features or too many parameters relative to the
 7 number of observations, or is forced to overextend its capabilities, i.e., make predictions with
 8 insufficient training data, the model runs the risk of overfitting (Srivastava et al., 2014). An
 9 overfitting model fits the training data “too well”, capturing noise and details that provide high
 10 accuracy on a training dataset, only to perform poorly once the model encounters “unseen” data in
 11 testing/validation (Xu and Liang, 2021). Scenarios where overfitting may be temporarily acceptable
 12 are those where: 1) model development is at its preliminary stages, where the interest is in a “proof of
 13 life” concept, 2) when the objective is to identify heavily-relied on features by the model, i.e., feature
 14 importance, or 3) in highly-controlled modeling environments where the expected data will be
 15 consistently similar to the training dataset. The latter is more likely in certain industrial applications
 16 and unlikely in the changing nature of hydrology.

17
 18 In contrast, underfitting occurs when a model is too simple to capture any patterns in the data, which
 19 can also lead to terrible performance in training, testing and validation. Underfitting can occur with
 20 inadequate model features, poor model complexity or when regularization techniques, (e.g., L1 or L2
 21 regularization), are over-used, making the model too rigid and unable to respond to changes in the
 22 data. Given the propensity of machine learning models to effectively learn the training data,
 23 underfitting is less of an issue in ML whereas overfitting can be widespread. In the following
 24 diagram, we present an example workflow to transition away from overfitting and towards
 25 generalizability. We further encourage modelers to actively transition towards making more
 26 generalizable models, which are in theory, more capable of performing well across diverse scenarios
 27 and datasets, which will become increasingly important with the persistence of climate extremes.

28
 29 **Response to ref #1, comment #1B:** We have added a few sentences (blue is new) to the Discussion
 30 subsection titled “ML as Knowledge Discovery” where we urge for TUURTs (Temporal, Unseen,
 31 Ungaged Region Tests):

32
 33 Our review finds that ML studies examining SWT have been conducted from a computational
 34 perspective, one with a focus on comparing techniques and performance metrics as opposed to
 35 explaining the nature of SWT dynamics or influencing processes. While it is understandable that not
 36 every ML-SWT paper aims to explain physical processes, we think the SWT community should come
 37 together and agree on a baseline of tests that all ML-SWT models should undergo for model
 38 robustness and transferability. Along these lines, we urge consideration of TUURTs (temporal,
 39 unseen, unged region tests) for future ML-SWT models as a helpful step towards not only better
 40 modeling practices but also increased model transparency and robustness. For this, we clarify that
 41 testing for “unseen” cases means testing only within the developmental dataset, whereas testing for
 42 “ungaged” cases means testing for new sites that have not been previously seen by the model at all.
 43 Recent ML-SWT studies have only applied one or two of the tests, but not all three (Topp et al.,
 44 2023; Hani et al., 2023, Souassi et al., 2023). Siegel et al. (2023), a non-ML SWT paper, tested for
 45 unged and unseen data but did not perform a temporal test. A relatively new study, Philippus et al.
 46 (2024), appears to be the only published SWT-ML study that purposefully applied TUURTs with
 47 some success.

48
 49
 50 *Disclaimer: I propose some additional literature (n = 4-5), and I am a coauthor on 1 of them. I do not*
 51 *view including that literature as mandatory, and only proposed additional sources based on their*
 52 *relevance to the content of this manuscript. I selected "No" to anonymity to avoid any appearance of*
 53 *subversive influence.*