

February 15, 2016

Dr. Dominic Mazvimavi
Department of Earth Sciences
Institute for Water Studies
University of the Western Cape
Bellville 7535
Republic of South Africa

Dear Dr. Mazvimavi,

We are submitting a revised manuscript for HESS-2015-413, entitled “Empirical streamflow simulation for water resource management in data-scarce seasonal watersheds” that addresses the comments provided by the two referees. Additionally, we’ve summarized how we responded to the specific comments in the text below, and attached a marked up version of the manuscript in track changes. We thank you for your time and consideration of our manuscript, and look forward to your decision.

Sincerely,

Julie Shortridge
Johns Hopkins University

Response to comments:

Reviewer 1:

1. *“The good performance of the climatological model is almost completely a result of the low interannual variability in the flow regime as evidenced by Figure 3, but this issue is never mentioned.”*

The reviewer makes a valid point that the original Figure 3 makes it appear as though there is very little interannual variability in the region’s hydrology. However, this actually is not the case, particularly when one considers the other rivers assessed as part of the study, each of which has approximately three times as much interannual variability (based on the coefficient of variation for annual flow volumes) as the Gilgel Abbay. The manuscript’s Figure 3 was just presented as an example to demonstrate the bias in standard formulation model predictions of wet season flows, and it could just as easily be replaced by a different river/time period to demonstrate that this phenomena. We have replaced this figure with a new one that does not suggest that interannual variability in the region is minimal. Additionally, we’ve edited Table 1 to incorporate information on the interannual variability of each river.

2. *“Some of the other comments in the paper about how the empirical models can be used to assess physical realism are also, in my opinion, rather tenuous. ‘Runoff increasing with*

higher precipitation levels and decreasing with higher temperatures' (page 19) is hardly a measure of physical realism... I therefore cannot agree with the authors that their models can be used to characterize 'watershed behaviour in a manner that could shed light on underlying physical processes' (page 18). If that is the case, what are the processes? Are the dry season processes groundwater driven or drainage from wetlands?"

We acknowledge that the models in their current formulation cannot be used to assess complex questions about physical hydrological process in the basins studied. However, we disagree that the simple relationships identified in the models are not measures of physical realism, and point towards similar evaluations that have been used to assess empirical model performance in the literature. For example, Han et al. (2007) explore how ANN flood forecasting models responds to a double-unit input of rain, finding that some formulations respond in a hydrologically meaningful way to increased rainfall intensity, while others do not. Similarly, Galelli and Castelletti (2013) describe how input variable importance can be used to highlight differences in hydrologic processes between an urbanized and forested watershed. In both cases, the relationships identified were fairly simplistic (e.g., higher runoff with greater rainfall intensity; shorter time of concentration in urban versus forested watersheds), but are still important steps in characterizing the mechanisms by which models make predictions so that they are not “black-boxes” and confirming that these mechanisms make physical sense. The text in section 4 has been updated to clarify this point.

3. *“The argument that these types of models are good for places where there are good climate data but poor physical data may be valid, but the real question is how often do such situations occur and if you have good flow data, why do you need a model to make water resources decisions.”*

While it is impossible to say exactly how often this situation occurs without a comprehensive review, this work was motivated by the specific data available in the Lake Tana region, where long records of flow and low-resolution climate data were available but detailed, ground-truthed spatial data on land cover and soil were not. In our experience this is not a rare situation: due to a combination of historical data centers (e.g., World Meteorological Organization reporting for climate and the Global Runoff Data Centre for streamflow) and more recent efforts to merge satellite data with *in situ* observations to monitor climate and hydrology (e.g., the Global Precipitation Climatology Project and the Global Land Data Assimilation System) one can often find acceptable climate data, even in data poor regions. Obtaining measurement-based estimates of soil hydraulic parameters or details on hydrologically-relevant land management activities can be more difficult. In this instance, there are two contexts in which historical flow data would be insufficient for decision making. In the short-term, models are needed to take advantage of seasonal climate forecasts to more efficiently manage hydropower and irrigation schemes. In the long-term, changes in land-cover and climate mean that historic data are unlikely to be representative of future flow conditions. Thus, any estimates of how proposed long-lived infrastructure will perform in coming decades requires models to translate climate and land cover conditions into flow. Text has been added to the introduction to clarify this point.

4. *“I am not sure about the value of the climate change scenarios as they appear to me to be very simplistic and add very little to the study.”*

We would like to clarify that these aren't climate change scenarios (which would describe plausible climate conditions expected to occur in the future), but instead measurements of the sensitivity and uncertainty of model predictions when forced with increasingly extreme climate data. Since one of the key motivations for using rainfall-runoff models is to understand how climate change may impact water resources, it is important to understand how model formulation contributes to this sensitivity and uncertainty. The analysis was thus kept intentionally simple, in an effort to avoid obscuring differences between models and implying that this analysis represented a projection of expected climate change impacts. While these issues could certainly be explored using actual downscaled climate model projections to make the assessment more representative of possible future impacts, we don't think that this additional complexity would add much to the inter-model comparison. The text of section 2.3 has been revised to clarify this point.

5. *“On page 6 the authors suggest that empirical models can provide more comprehensive uncertainty analysis results. Why when there are many recent examples of rainfall runoff models being used for uncertainty analysis and the assessment of model results from a behavioural and non-behavioural standpoint.”*

We acknowledge the reviewer's point that there have been a number of instances where uncertainty assessment has been conducted using physical rainfall-runoff models, but the statement on page 6 was referring specifically to the Lake Tana basin where there has been relatively little assessment of uncertainty in hydrologic modeling studies. The text in section 2.1 has been revised to clarify this point.

6. *“There are many places in the text where the word 'data' is treated as singular, while it should always be treated a plural (i.e. 'these data', 'date were', 'data area', etc.).”*

We thank the reviewer for identifying this error; it has been corrected in the revised manuscript.

7. *“The reference to the estimates of rainfall intensity on page 7 should be removed as this method will never give a proper estimate of intensity.”*

We agree with the reviewer that the method is a very rough approximation of actual rainfall intensity. However, when presenting this work we have received multiple questions about whether rainfall intensity was considered in the evaluation, and thus think that this should remain in the manuscript to demonstrate that we did consider intensity to the degree that the available data allowed and found that it was not a useful addition to the models.

8. *“Page 8 refers to a log transformation of monthly streamflow to get a better match to normal, however, the distribution properties of the monthly flow data are not assessed.”*

This text has been revised.

9. *“If NSE is considered such a bad statistic, why not use something else. Even NSE based on log transformed values can remove some of the bias to high wet season flows.”*

We used NSE based on raw flow values (rather than log-transformed) because that is what has been used in other modeling studies conducted in the basin and thus seemed like the most appropriate metric for comparison. It should be noted that we don't necessarily disagree with the use of NSE as a metric, but rather the assumption that an NSE score greater than 0.5 indicates good model performance. Text regarding the choice of error metrics has been added to section 2.3.

Reviewer 2:

General comments:

1. *The description of the data-driven models (line 14, page 11091 – line 20, page 11092) is too synthetic and thus prevents the reader from understanding the experimental set-up (e.g., Table 2) as well as some of the results reported in Section 3...The experimental set-up is described only partially and some of the adopted techniques require more parameters than those listed in Table 2.*

To ensure that the work is clear and reproducible, we have expanded Table 2 so that it includes all of the parameters used to fit the models (rather than just those that were optimized through cross validation). The description of modeling approaches on pages 11091-11092 is meant to only provide the reader with a brief background on the methods used and provide references where each approach is discussed in more detail. The references cited in this section provide details on how each model is trained; this information, provided with the details on the specific parameterizations used in Table 2, should be sufficient for understanding model development process.

2. *I have some doubts regarding the second formulation (Equations 2-3). Streamflow anomalies are calculated by (a) subtracting the long-term average streamflow and (b) dividing this number by the long-term standard deviation. However, the streamflow process appears to be non-stationary... the changes in land use have an impact on the rainfall-runoff process, while the long-term average and standard deviation are calculated on the hypothesis of a stationary process. I think that the authors should elaborate on this point.*

The referee is correct that streamflow processes in the region are most likely non-stationary due to changes in land cover over decadal time scales as well as the influence of rising temperatures. However, these changing conditions are incorporated into the calculation of the streamflow anomaly value itself, since this value is a function of temperature, rainfall, and agricultural land cover. Thus, while the conversion from streamflow anomaly to raw streamflow value in CMS uses stationary measurements of long-term average and standard deviation, the calculation of the anomaly value itself does not rely on any assumption about stationary conditions. Text has been added to Section 2.2 to clarify this point.

3. *I do not understand why they have not used an additional (and better) metric, such as KGE.*

The use of NSE was included because it is the most widely used error metric in modeling studies conducted in the region, and provided a rough point of comparison between these models and physical models that had been previously developed for the region. MAE was included as an error metric because it provides a simple and easily interpretable measure of error on the same scale as observed flow volumes. The use of alternative error metrics has been discussed extensively in the literature (for instance Pushpalatha et al., 2012; Mathevet et al., 2006; Criss and Winston, 2008). While model evaluation in terms of alternative or additional error metrics could provide interesting insights into what is contributing to predictive capabilities of different model formulations, the objective of this paper is to look

beyond predictive capability and instead compare model formulations in terms of error structure and uncertainty. Examination of the KGE performance metric (Gupta et al., 2009), for example, confirms that models outperform climatology in all watersheds, though the specific ranking of model performance does change in some cases. Text has been added to section 2.3 that elaborates on this point.

4. *Why does the climatology model perform so well? Given the results reported in Table 3, one might conclude that complex data-driven models are not needed since a simple climatological model can get excellent values of NSE and MAE.*

The climatology model does well because seasonality accounts for such a large portion of the variability in monthly flow, a phenomenon discussed by Legates and McCabe (1999) and Schaefli and Gupta (2007). However, this model does not account for any degree of interannual variability nor the possibility for non-stationary conditions caused by changing land cover and climate, and thus is unsuitable for streamflow simulation over the short term (eg., based on seasonal climate forecasts) or long term (due to land cover and climate change). Text has been added regarding this point to section 3.1 and section 4.

Specific Comments:

1. *The title does not fully represent the content of the paper; in particular, water management issues are not explored in the study.*

We respectfully disagree with the reviewer on this issue, as the issues of model interpretability, bias, and uncertainty are very important in terms determining whether a model is fit-for-purpose. For this reason, we have left the title unchanged.

2. *Line 2, page 11084. Can you give an example of the “certain methods” mentioned here?*
Text revised to say “machine-learning methods”.

3. *Line 4, page 11084. “Data” should be used as the plural form of ‘datum’.*
Corrected.

4. *Line 4-10, page, 11084. I do not completely agree with this statement. There is an extensive body of literature on the application of data-driven techniques to streamflow modelling problems. See, for example, Elshorbagy et al. (2010). Whilst model interpretability and uncertainty have received somewhat less attention, there have been studies focusing on such aspects; see, for example, Wilby et al. (2003) or Taormina and Chau (2015).*

While we acknowledge that there are studies evaluating other methodologies and aspects of model performance beyond predictive accuracy, we stand by our statement that the majority of papers that use machine learning for streamflow simulation apply ANNs and focus on predictive accuracy. For instance, a review by Solomatine and Ostfield (2008) of data-driven hydrologic modeling states that multilayer perceptron ANNs “have become the most popular machine learning tool... and are known to have several dozens of successful applications in river basin management and related problems.” Additionally, while hundreds of manuscripts on ANNs have been published (see review by Maier et al., 2010), the number that do consider interpretability and uncertainty have been a relatively small

percentage of these. While a full discussion of these issues and referencing of supporting literature is outside the scope of the abstract, additional text supporting this statement has been added to the introduction.

5. *Line 11, page 11084. I think the authors should explicitly mention the “machine learning” techniques used in their study.*
Added.
6. *Line 20, page 11084. This sentence may be misleading, since the study does not carry out a climate impact assessment.*
Revised.
7. *Line 4, page 11085. Are the authors referring to Genetic Programming (Babovic et al., 2005)? Genetic algorithms are heuristic optimization techniques; as such, they cannot be directly employed for data-driven streamflow predictions.*
This has been corrected.
8. *Line 24-27, page 11085. Yes, but this why there exists a variety of techniques (cross-validation, bootstrapping etc.) aimed at minimizing/reducing overfitting problems.*
Text revised to mention this.
9. *Line 21-22, page 11086. Yes, this why different (or multiple) objective functions should be considered when training a model (De Vos and Rientjes, 2008).*
Agreed. However, the use of multiple objective functions cannot capture every aspect of model performance that may be relevant for a given planning purpose, such as interpretability or uncertainty. Text has been added expanding on this point.
10. *Line 1-3, page 11087. Again, I believe that during the past 5-10 years, several studies on data-driven streamflow forecasts not only focused on improving model performance, but also on improving our understanding of the models structure, thus supporting the interpretation of the underlying physical processes.*
The text has been updated to reference some of the studies that focus on model structure and its interpretation. However, one limitation with these studies is that they generally focus on a single method, and thus provide little insight into the relative ease with which the structure of different model types can be investigated. Text has been added discussing this point.
11. *Line 7-9, page 11087. Can the authors expand the literature review on the use of data-driven techniques on non-temperate regions?*
Additional text on previous work done in these regions has been added.
12. *Line 16, page 11087. Can the authors provide more details on “relevant landscape change”?*
Added.
13. *Line 1, page 11088. This section should be named ‘Data and Methods’.*
Revised.

14. *Line 20-21, page 11088. Can you give an example of these infrastructures?*
Added.
15. *The first part of Section 2.2 is about data, not models. Why not splitting it into two sections focusing on data and models, respectively?*
Since only the first paragraph of this section discusses data, we do not think it makes sense to split it into a separate section, but we have changed the title of this section to “Data and Model Development”
16. *Line 22, page 11089. “monthly daily average temperature”?*
Revised to say “monthly average temperature”
17. *Line 28, page 11089. It should be “these data”.*
Revised
18. *Line 4, page 11090. Does the land cover vary on an annual basis?*
Yes – text has been revised to clarify this.
19. *Line 15-16, page 11090. This comment is about model results, not data. It should be moved to the results section.*
This is technically correct, but is presented more as a justification for using admittedly limited land cover data in the analysis, rather than a main finding of the investigation. For this reason, we think it makes more sense in the methodology section and have left it there.
20. *At the end of the first paragraph (Section 2.2) authors should clearly state what the number of available observations (for each catchment) is.*
Added.
21. *Moreover, do the authors have an estimate of the time of concentration (of each catchment)? This relates to the time lags adopted for the precipitation in model (1) and (2)-(3).*
The time of concentration for each catchment is not known, but the declining influence of lagged climatic variables at two months prior indicate that climate conditions from beyond this time period are unlikely to contribute to flow variability. Text has been added clarifying this point.
22. *Line 27-28, page 11090. Do the streamflow data follow a log-normal distribution?*
This text has been revised – while the distribution of streamflow more closely resembles a log-normal distribution when compared to a normal distribution, the presence of high values makes the distribution somewhat asymmetric. However, this transformation was still necessary to avoid predicting negative flow values.
23. *Line 13, page 11091. It should be “Six”, not “Seven”.*
Revised

24. *Line 6, page 11093. It should be Table 2, not Table 1. This error is repeated through-out the manuscript.*
This has been corrected.
25. *Line 18, page 11093. What is the reason for adopting the MAE?*
MAE was included as an error metric because it provides a simple and easily interpretable measure of error on the same scale as observed flow volumes. Text has been added to clarify this point.
26. *Line 20-22. Why?*
Because the high degree to which seasonality contributes to variability in flow means that it is easy for relatively uninformative models to obtain high NSE values as long as they are able to capture basic seasonality. Text has been added to clarify this.
27. *Line 19, page 11095. What is the “delta-change method”? I think that a short explanation is needed.*
Added.
28. *Line 25-28, page 11097. One of the most common mechanisms for understanding the importance and influence of covariates is input variable selection - see, Wu et al. (2014) and Galelli et al. (2014).*
Agreed. However, because this section is on the results of the covariate influence assessment (which did not include variable selection since that is not an established methodology for all models evaluated), additional text on this point has been added to the introduction.
29. *Line 27-28, page 11099. Is there any reason behind this?*
This would mean that the model predictions are largely based on land cover and precipitation, which would suggest that these variables are the dominant controls on streamflow variability
30. *Line 6-7, page 11104. This comment is not necessary.*
We disagree that the comment is not necessary, as the need to write special code for interpreting machine-learning models necessarily requires additional skill and background in their implementation that may not be common amongst practitioners. Thus, the ability for this information to be obtained in a relatively straightforward way is another aspect of a model that makes it fit-for-purpose.

References

- Babovic, Vladan, and Maarten Keijzer. "Rainfall Runoff Modeling Based on Genetic Programming." *Encyclopedia of Hydrological Sciences* (2005).
- Criss, R.E., Winston, W.E., 2008. Do Nash values have value? Discussion and alternate proposals. *Hydrological Processes* 22, 2723–2725.
- De Vos, N. J., and T. H. M. Rientjes. "Multiobjective training of artificial neural networks for rainfall-runoff modeling" *Water resources research* 44.8 (2008).
- Elshorbagy, A., et al. "Experimental investigation of the predictive capabilities of data driven modeling techniques in hydrology-Part 1: Concepts and methodology." *Hydrology and Earth System Sciences* 14.10 (2010): 1931-1941.
- Galelli, Stefano, et al. "An evaluation framework for input variable selection algorithms for environmental data-driven models." *Environmental Modelling & Software* 62 (2014): 33-51.
- Gupta, H.V., et al., 2009. Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. *Journal of Hydrology*, 377(1), pp.80-91.
- Legates, D. R. and McCabe Jr, G. J.: Evaluating the use of “goodness-of-fit” measures in hydrologic and hydroclimatic model validation, *Water Resour. Res.*, 35, 233–241, 1999.
- Mathevet, T., et al. 2006. A bounded version of the Nash–Sutcliffe criterion for better model assessment on large sets of basins. In: Andréassian, V., et al. (Eds.), *Large Sample Basin Experiment for Hydrological Model Parameterization: Results of the Model Parameter Experiment – MOPEX*. IAHS Publ, p. 567.
- Pushpalatha, R., et al. 2012. A review of efficiency criteria suitable for evaluating low-flow simulations." *Journal of Hydrology* 420: 171-182.
- Schaefli, B. and Gupta, H. V.: Do Nash values have value?, *Hydrol. Process.*, 21, 2075–2080, doi:10.1002/hyp.6825, 2007.
- Taormina, Riccardo, and Kwok-Wing Chau. "ANN-based interval forecasting of streamflow discharges using the LUBE method and MOFIPS." *Engineering Applications of Artificial Intelligence* 45 (2015): 429-440.
- Wilby, R. L., R. J. Abrahart, and C. W. Dawson. "Detection of conceptual model rainfall runoff processes inside an artificial neural network." *Hydrological Sciences Journal* 48.2 (2003): 163-181.
- Wu, Wenyan, Graeme C. Dandy, and Holger R. Maier. "Protocol for developing ANN models and its application to the assessment of the quality of the ANN model development process in drinking water quality modelling." *Environmental Modelling & Software* 54 (2014): 108-127.

1 **Empirical streamflow simulation for water resource**
2 **management in data-scarce seasonal watersheds**

3

4 **J. E. Shortridge,¹ S. D. Guikema,² and B. F. Zaitchik³**

5 [1]{Department of Geography and Environmental Engineering, Johns Hopkins University,
6 Baltimore, USA}

7 [2]{Department of Industrial and Operations Engineering, University of Michigan, Ann
8 Arbor, USA}

9 [3]{Department of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, USA}

10 Correspondence to: J. E. Shortridge (jshortridge@jhu.edu)

11

12

1 Abstract

2 In the past decade, ~~machine-learning~~ertain methods for empirical rainfall-runoff modeling
3 have seen extensive development and been proposed as a useful complement to physical
4 hydrologic models, particularly in basins where data to support process-based models are
5 limited. However, the majority of research has focused on a small number of methods, such as
6 artificial neural networks, despite the development of multiple other approaches for non-
7 parametric regression in recent years. Furthermore, this work has ~~generally~~often evaluated
8 model performance based on predictive accuracy alone, while not considering broader
9 objectives such as model interpretability and uncertainty that are important if such methods
10 are to be used for planning and management decisions. In this paper, we use multiple
11 regression and machine-learning approaches (including generalized additive models,
12 multivariate adaptive regression splines, artificial neural networks, random forests, and M5
13 ubist models) to simulate monthly streamflow in five highly-seasonal rivers in the highlands
14 of Ethiopia and compare their performance in terms of predictive accuracy, error structure and
15 bias, model interpretability, and uncertainty when faced with extreme climate conditions.
16 While the relative predictive performance of models differed across basins, data-driven
17 approaches were able to achieve reduced errors when compared to physical models developed
18 for the region. Methods such as random forests and generalized additive models may have
19 advantages in terms of visualization and interpretation of model structure, which can be useful
20 in providing insights into physical watershed function. However, the uncertainty associated
21 with model predictions under ~~climate change~~extreme climate conditions should be carefully
22 evaluated, since certain models (especially generalized additive models and multivariate
23 adaptive regression splines) becoame highly variable when faced with high temperatures.

24

25

1 1 Introduction

2 Hydrologists and water managers have made use of observed relationships between
3 rainfall and runoff to predict streamflow ever since the creation of the rational method in the
4 19th century (Beven, 2011). However, the development of increasingly sophisticated machine
5 learning techniques, combined with rapid increases in computational ability, has prompted
6 extensive research into advanced methods for data-driven streamflow prediction in the past
7 decade. Artificial neural networks (ANNs), regression trees, and support vector
8 machines~~genetic algorithms~~ have been shown to be powerful tools for predictive modeling
9 and exploratory data analysis, particularly in systems that exhibit complex, non-linear
10 behavior (Solomatine and Ostfield, 2008; Abrahard and See, 2007).

11 While distributed physical models that accurately represent hydrologic processes can still
12 be considered the gold standard for rainfall runoff modeling, empirical models can be a useful
13 tool in contexts where there is limited data on physical watershed processes but long time-
14 series of precipitation and streamflow (Iorgulescu and Beven, 2004). The development of
15 historical data centers and more recent efforts to merge satellite data with *in situ* observations
16 to monitor climate and hydrology has made acceptable climate and streamflow data more
17 widely available in data poor regions. Because obtaining measurement-based estimates of soil
18 hydraulic parameters or details on hydrologically-relevant land management activities can be
19 more difficult, empirical models may be particularly useful in these locations. While many
20 criticize these approaches as “black boxes” with no relationship to underlying physical
21 processes (See et al., 2007), a number of studies have demonstrated how empirical approaches
22 can be used to gain insights about physical system function (e.g., Han et al., 2007; Galelli and
23 Castelletti, 2013a). Additionally, improvements in interpretation and visualization methods
24 can make complex models more easily interpretable (Sudheer and Jain, 2004; Jain et al.,
25 2004). Finally, data-driven models can be useful in identifying situations where observed data
26 disagree with what would be predicted based on conceptual models, and thus identify
27 assumptions regarding runoff generation processes that may be incorrect (Beven 2011).

28 While there have been some applications of alternative machine learning methods, such
29 as support vector machines (Asefa et al., 2006; Lin et al., 2006) and regression-tree based
30 approaches (Iorgulescu and Beven, 2004; Galelli and Castelletti, 2013a) for streamflow
31 simulation, the vast majority of research has focused on artificial neural networks (Solomatine
32 and Ostfield, 2008). While they have demonstrated impressive predictive accuracy in a

1 number of different contexts, excessive parameterization of ANNs can result in overfit
2 models that are not generalizable to unseen data (Iorgulescu and Beven, 2004; Gaume and
3 Gosset, 2003). While methods exist to avoid overfitting, such as cross validation and
4 bootstrapping, these methods are not always employed (Solomatine and Ostfield, 2008). A
5 review by Maier et al. (2010) found that relatively few studies evaluated model performance
6 based on parameters such as Akaike information criterion that would lead to parsimonious
7 models that are likely to be more generalizable and interpretable. This can lead to complex
8 models that only result in modest improvements (or no improvements at all) over much
9 simpler approaches (Gaume and Gosset, 2003; Han et al., 2007).

10 Even outside of a hydrology context, it has been argued that ANNs are better suited for
11 problems aimed at prediction without any need for model interpretation, rather than those
12 where understanding the process generating predictions and the role of input variables is
13 important (Hastie et al., 2009). Given the importance that this interpretation plays in
14 understanding the contexts in which a hydrologic model is appropriate and reliable, the strong
15 opinions surrounding the use of ANNs for water resources management are perhaps not
16 surprising. To address this issue, a number of studies have focused on highlighting the
17 structure and mechanism by which machine learning models make predictions to confirm
18 their physical realism and gain insight into physical watershed function. For example, some
19 studies have demonstrated how internal ANN structure corresponds to physical hydrologic
20 processes (Wilby et al. 2003; Jain et al., 2004; Sudheer and Jain, 2004), while others have
21 shown how variable selection and importance can be used to gain insights about model
22 structure and runoff generating processes (Galelli and Castelletti, 2013a and 2013b). While
23 these studies demonstrate that a number of methods exist for characterizing model structure,
24 they generally focus on a single model type and thus provide little insight into the
25 comparative ease with which different model types can be interpreted.

26
27 While a number of comparison studies exist that apply multiple empirical models to a
28 given problem, finding generalizable insights from these studies is hindered because ~~they~~ of
29 the limited number of models and datasets evaluated. Perhaps the most comprehensive
30 comparison to date is that of Elshorbagy et al. (2010a and 2010b), who compared six methods
31 for data-driven modeling of daily discharge in the Ourthe River in Belgium. This work found
32 that linear models were able to perform comparably to much more complex methods when the

1 data content of the models were limited, or when system input-output behavior was close to
2 linear. However, other studies have demonstrated the value of using more complex
3 approaches when modeling more complex rainfall-runoff behavior (e.g., Abrahart and See,
4 2007; Asefa et al., 2006). The differing results obtained across these studies indicate that no
5 single method is likely to be suitable for all basins, timescales, or applications.

6 However, it is important to recognize that predictive accuracy alone is not necessarily
7 sufficient justification for applying a model to a given problem. Models should not only be
8 accurate, but also be fit-for-purpose (Beven, 2011; Van Griensven et al., 2012). For instance,
9 accurate representation of low return period flows is more important in a flood forecasting
10 model than one aimed at predicting average amounts of water available for withdrawal and
11 human consumption. Similarly, the ability to provide insights into physical watershed
12 function may be more important in basins where land-use change could alter the hydrologic
13 regime, compared to a basin that is heavily urbanized and expected to remain so. The use of
14 multiple objective functions in training data-driven models can address this to some degree by
15 identifying models that provide sufficient balance between different performance objectives,
16 such as accurate representation of different portions of the flow hydrograph (De Vos and
17 Rientjes, 2008). However, more refined model training procedures will not necessarily
18 address other aspects of model performance that make it suitable for planning purposes, such
19 as interpretability (Solomatine and Ostfield, 2008). This is particularly true for data-driven
20 models; as pointed out by Solomatine and Ostfield (2008), there is a need for additional
21 research on making models more understandable and useful for water managers. More
22 comprehensive consideration of model strengths and limitations should be standard practice
23 in model development and selection, rather than simply evaluating global error metrics.

24 In this work, we compare six methods for empirical streamflow prediction (linear models,
25 generalized additive models, multivariate adaptive regression splines, random forests, M5
26 model trees and ANNs) in their ability to predict monthly streamflow in five rivers in the
27 Lake Tana basin in Ethiopia. This study region was selected as it provides insights into the
28 use of data-driven models for streamflow simulation in tropical regions of the world that are
29 underrepresented in existing studies; for instance, a review of 210 articles on water resource
30 applications of ANNs found that over three quarters of the studies evaluated were conducted
31 in North America, Europe, Australia, or temperate East Asia (Maier et al., 2010). Existing
32 studies conducted in tropical regions generally apply a single methodology to the basin of

1 interest and evaluate predictive accuracy alone (see for instance, Machado et al., 2011;
2 Chibanga et al., 2003; Antar et al., 2006; Aqil et al., 2007), making it difficult to find
3 generalizable insights into the relative advantages of different modelling approaches in these
4 regions. to provide a counterpart to previous comparative studies that have largely focused on
5 rivers in temperate regions. Better development of data-driven models for these regions has
6 the potential to be particularly valuable because data limitations and complex hydrodynamic
7 processes often hinder the use of physical watershed models, but Furthermore, physically-
8 based hydrological process models are a challenge in these basins due to data limitations—
9 soil and vegetation parameters are poorly characterized and high frequency, spatially-
10 distributed precipitation estimates are highly uncertain—and complex hydrodynamic
11 processes, including lake backwater effects, that are neglected by most watershed models.
12 There are, however, relatively long time series of streamflow, available, and estimates of
13 historical precipitation and temperature may be available at a monthly timescale. These
14 data, combined with information on relevant landscape change (in particular, the expansion of
15 agricultural land cover), can be leveraged to create reasonably accurate empirical models.

16 Models are compared not only in terms of their predictive accuracy, but also in terms of
17 model error structure and the implications that this structure may have for water resource
18 applications. Additionally, we evaluate the methods by which model structure and predictor
19 variable influence can be evaluated to gain insights into physical system function for each
20 model type. Finally, we assess the suitability of using different model types for climate
21 change impact assessment by comparing model uncertainty in projections made for
22 increasingly extreme climate conditions. The overall objective of this research is not to
23 identify a single “best” model, but rather to highlight some of the strengths and limitations of
24 different approaches, as well as demonstrate important issues that should be kept in mind for
25 model comparisons in the future

26 **2 Data and Methods**

27 **2.1 Study Area**

28 Lake Tana is located at an elevation of approximately 1800 meters in the highlands of
29 northwest Ethiopia (Fig. 1). The catchment draining to the lake encompasses approximately
30 12,000 square kilometers, and the four main tributaries providing water to the lake are the
31 Gilgel Abbay (including its tributary, the Koga River), Ribb, Gumara, and Megech Rivers.

1 Collectively, these rivers account for 93% of the inflow to the lake (Alemayehu et al., 2010).
2 Ninety percent of rainfall in the basin occurs during the wet season from May until October,
3 and there is significant interannual variability in precipitation with annual rainfall levels
4 ranging from below 1000 mm to over 1800 mm (Achenef et al., 2013). Population growth and
5 expansion of agricultural and pastoral land use in the region has resulted in substantial
6 deforestation and land degradation, with agricultural, pastoral and settled land cover
7 comprising over 70% of the basin's surface area (Rientjes et al., 2011; Garede and Minale,
8 2014; Gebrehiwot et al., 2010). There is some evidence that this has impacted the hydrology
9 of the rivers draining into the lake (Gebrehiwot et al., 2010). A summary of basin
10 characteristics for the evaluation period of 1960-2004 is presented in Table 1.

11 Approximately 2.6 million people live in the basin, and are largely settled in rural
12 areas and reliant on rainfed subsistence agriculture. This makes the region quite vulnerable to
13 climate variability and change, and a number of water resources infrastructure projects are
14 planned to better manage this vulnerability and support economic development (Alemayehu
15 et al., 2010). This includes the recent construction of the Tana-Beles hydropower transfer
16 tunnel and the Koga River irrigation reservoir, as well as five other reservoirs planned for
17 construction in the next 10 to 20 years (Alemayehu et al., 2010). To better understand the
18 potential implications of this development, extensive effort has been put towards developing
19 rainfall-runoff models for the Lake Tana basin, as well as other areas of the Ethiopian
20 highlands with similar characteristics (van Griensven et al., 2012). Many of these studies rely
21 on Soil and Water Assessment Tool (SWAT) models, although there are some that use water
22 balance approaches (Van Griensven et al., 2012). While these models have in some cases
23 demonstrated reasonably high accuracy, previous evaluations were largely based on Nash-
24 Sutcliffe Efficiency (NSE; Nash and Sutcliffe, 1970) which can be a flawed performance
25 metric in highly seasonal watersheds (Schaeffli and Gupta, 2007; Legates and McCabe, 1999).
26 More importantly, the limited data available for physical parameterization of these models
27 required a heavy reliance on model calibration, which sometimes resulted in parameterization
28 schemes that are inconsistent with physical understanding of the region's hydrology
29 (Steenhuis et al., 2009; van Griensven et al., 2012). Furthermore, a number of studies relied
30 on empirical relationships such as curve numbers and the Hargreaves equation that were
31 developed for temperate regions (e.g., Mekonnen et al., 2009; Setegne et al., 2009). While
32 these limitations are likely to introduce considerable uncertainty into model projections,
33 particularly in situations where climatic or environmental conditions differ from those

1 experienced in the calibration period, few studies from this region of Ethiopia include any sort
2 of uncertainty analysis in model predictions. Empirical models could provide a useful
3 complement to physical models developed for the region by providing insights into physical
4 system function and allowing for more comprehensive uncertainty analysis.

5 **2.2 Data and Model Development**

6 Models were developed using monthly streamflow, climate, and land cover data for
7 the period from 1961 to 2004, resulting in 528 monthly observations. In each of the five major
8 rivers in the basin, we developed empirical models that estimated monthly streamflow as a
9 function of climate conditions and agricultural land cover in each basin. Monthly streamflow
10 data ~~were~~ taken from historic stream gauge records for each basin, as reported in
11 feasibility studies developed for proposed irrigation projects (Alemayehu, 2010). Historic data
12 for monthly ~~daily~~ average temperature, monthly total precipitation, and monthly wet days in
13 each river basin were derived from the University of East Anglia Climate Research Unit
14 (CRU) TS3.10 gridded meteorological fields (Harris et al., 2014), which are based on
15 meteorological station observations. Historic estimates of rainfall intensity were also
16 calculated by dividing monthly total precipitation by CRU TS3.10 records of the number of
17 wet days in that month, ~~but was. However, this data was~~ found to be highly correlated with
18 monthly precipitation and did not result in significant improvements to the predictive
19 accuracy of tested models. ~~Thus, it was, and was thus~~ not included in the final model
20 formulations. Finally, to account for historic increases in agricultural and pastoral land cover
21 that have occurred in the basin, the percentage of land cover used for any crop or grazing was
22 estimated from historic land cover analyses described by Rientjes et al. (2011), Gebrehiwot et
23 al. (2010), and Garede and Minale (2014). These studies used historic aerial photos and
24 satellite images to estimate land cover changes in the Ribb, Gilgel Abbay, and Koga basins
25 from the periods of 1957 to 2011. The percentage of agricultural land cover was interpolated
26 for years when data ~~were~~ ~~wasn't~~ available, and the value of agricultural land cover in the
27 two basins without data was assumed to be equal to average agricultural land cover in the
28 basins with data. Land cover was assumed to change on an annual, rather than monthly basis.
29 While this approach is prone to errors that could stem from differing rates of land use change
30 through time and between basins, it does provide a mechanism for capturing the long-term
31 trend of expanding agricultural land cover that has been observed throughout the Ethiopian
32 highlands when detailed land-cover data ~~are~~ ~~is~~ unavailable. Including this data improved out-

1 of-sample predictive accuracy of the models, further suggesting that it was a valuable
2 addition.

3 Two general formulations for the empirical models were evaluated. The first (referred
4 to below as the standard model formulation) was

$$5 \quad \log(Q_{b,t}) = f(P_{b,t}, P_{b,t-1}, P_{b,t-2}, T_{b,t}, T_{b,t-1}, T_{b,t-2}, AgLC_{b,t}) + \varepsilon_{b,t} \quad (1)$$

6 where $Q_{b,t}$ is the monthly streamflow in river b at time period t , $P_{b,t}$ and $T_{b,t}$ are the monthly
7 total precipitation and average temperature in river basin b at time period t , $AgLC_{b,t}$ is the total
8 percentage of agricultural land cover in basin b at time t , and $\varepsilon_{b,t}$ is the model error. The
9 subscripts $t-1$ and $t-2$ indicate lagged measurements from one and two months prior, and were
10 included to roughly account for storage times longer than one month that could impact
11 streamflow in each river. While the exact time of concentration is not known in each basin,
12 the minor influence of of climate conditions at two months prior suggest that climate
13 conditions from beyond this time period do not contribute significantly to flow variability.

14 The function f represents a general function that differed between the specific models assessed
15 and is discussed in more detail below. The logarithm of monthly streamflow was used as a
16 response variable to keep model predictions positive, ~~and make the response variable better~~
17 ~~match a normal distribution, a requirement for the use of Gaussian linear models and~~
18 ~~generalized additive models.~~

19 In the second formulation, streamflow and climate anomalies were used as the
20 response and predictor variables to better account for the highly seasonal nature of streamflow
21 and precipitation in the region. Streamflow anomalies were calculated for each observation by
22 subtracting the long-term average streamflow for that month (m) from the observed value and
23 dividing this number by the long-term standard deviation of that month's streamflow as in Eq.
24 (2). This procedure was repeated for precipitation and temperature, and these values were then
25 used to fit models of the form described in Eq. (3). It should be noted that although this
26 formulation uses long-term averages and standard deviations to convert anomaly values to
27 flow volumes, the anomaly values themselves are calculated based on climatic and land cover
28 conditions that are nonstationary through time.

$$29 \quad Q_{b,t}^{AN} = \frac{Q_{b,t} - \bar{Q}_{b,m}}{sd(Q_{b,m})} \quad (2)$$

30

$$Q_{b,t}^{AN} = f(P_{b,t}^{AN}, P_{b,t-1}^{AN}, P_{b,t-2}^{AN}, T_{b,t}^{AN}, T_{b,t-1}^{AN}, T_{b,t-2}^{AN} AgLC_{b,t}) + \varepsilon_{b,t} \quad (3)$$

Sixteen different types of models were compared using each formulation in each basin:

1. A Gaussian linear regression model (GLM) using the basic stats package in the R statistical computing software (R Development Core Team, 2014)
2. Gaussian generalized additive model (GAM): GAMs are a semi-parametric regression approach where the response variable is estimated as the sum of smoothing functions applied over predictor variables. These functions allow the model to capture non-linear relationships between the predictor and response variables without *a priori* assumptions about the form (eg., quadratic, logarithmic) of these functions, and are fit using penalized likelihood maximization to prevent model overfitting (Hastie and Tibshirani, 1990). GAMs were fit using the mgcv package in R (Wood, 2011).
3. Multivariate adaptive regression splines (MARS): MARS are a non-parametric regression approach where the response variable is estimated as the sum of basis functions fit to recursively partitioned segments of the data (Friedman, 1991). MARS models were fit using the earth package in R (Milborrow, 2015).
4. Artificial neural network (ANN): ANNs are a non-parametric regression approach represented by a network of nodes and links that connects predictor variables to the response variable. Each link in the network represents a function that maps the input nodes into the output node (Ripley, 1996). ANN models were fit using the nnet package in R (Venables and Ripley, 2013).
5. Random forest (RF): Random forests are a rule-based, non-parametric regression approach where the model prediction is created by averaging the predicted value from multiple regression trees which are trained on separate bootstrapped resamples of the data. Each tree is fit using a small, randomly selected subset of predictor variables, resulting in reduced correlation between trees (Breiman, 2001). Random forest models were fit using the randomForest package in R (Liaw and Wiener, 2002).
6. M5 model: M5 models are a rule-based, non-parametric regression approach that fits a linear regression model to each terminal node of a regression tree (Quinlan, 1992). M5 models were fit using the Cubist package in R (Kuhn et al., 2014).

1 7. Climatology model: A climatology model that simply predicted each month's
2 streamflow as equivalent to the long-term average streamflow for that month was
3 included for comparison purposes.

4 **2.3 Model Evaluation**

5 When using non-parametric regression approaches, it is important to avoid overfitting a
6 model to a given dataset because this can result in large errors in out-of-sample predictions
7 (Hastie et al., 2009). To avoid model overfit, the caret package in R (Kuhn, 2015) was used to
8 determine model parameters for the MARS, ANN, RF and M5 models. This package uses
9 resampling to evaluate the effect that model parameters have on the model's predictive
10 performance and chooses the set of parameters that minimizes out-of-sample error (Kuhn
11 2015). In this evaluation, 25 bootstrap resamples of the training dataset were generated for
12 each parameter value to be assessed. A model was fit using each bootstrap sample and used to
13 predict the remaining observations, and the parameter values that minimized average RMSE
14 across all resamples. Details on the specific parameters evaluated for each model are
15 presented in Table 42. While the development of more complex structures are possible for
16 some models, this process can result in over-parameterization and poor model performance
17 (Gaume and Gosset, 2003; Han et al., 2007). Additionally, the use of a standardized
18 parameterization procedure allows for a more even comparison between different model
19 types.

20 The predictive ability of each model was assessed using 50 random holdout cross-
21 validation samples. In each sample, a random selection of years were selected, and
22 observations from these years were removed ("held-out") from the dataset. The size of the
23 held-out sample ranged from 1 to 9 years. Each model was then fit to the remaining portion of
24 the data, using the caret package described above to determine model parameters for the
25 MARS, ANN, RF and M5 models. These models were then used to predict streamflow for the
26 held-out portion of the data, and both the mean absolute error (MAE) and NSE were
27 calculated after transforming model predictions after back to the original streamflow units.

28 ~~While NSE values are acknowledged to be a flawed performance metric in highly seasonal~~
29 ~~watersheds (Schaefer and Gupta, 2007; Legates and McCabe, 1999), this metric was included~~
30 ~~to provide a rough comparison of how empirical model performance compared to the~~
31 ~~performance of physical models developed for the region.~~ Mean MAE and NSE were
32 calculated for each model across the 50 cross-validation samples and used to choose the

1 model with the highest predictive accuracy in each basin. This cross-validation procedure
2 provides a mechanism for evaluating how well a model will generalize to an unseen set of
3 data while avoiding some of the problems that can arise from the use of a single calibration
4 and validation dataset (Elshorbagy et al., 2010a; Han et al., 2007).

5 MAE was included as an error metric because it provides a simple and easily
6 interpretable measure of error on the same scale as observed flow volumes. While NSE values
7 are acknowledged to be a flawed performance metric in highly seasonal watersheds where
8 seasonal fluctuations contribute to a substantial portion of flow variability (Schaefli and
9 Gupta, 2007; Legates and McCabe, 1999), this metric was included to provide a rough
10 comparison of how empirical model performance compared to the performance of physical
11 models developed for the region. The use of alternative error metrics has been discussed
12 extensively in the literature (for instance Pushpalatha et al., 2012; Mathevet et al., 2006; Criss
13 and Winston, 2008), and could provide additional insights into what contributes to predicitive
14 capabilities of different model formulations. However, this work examined predicitive
15 accuracy based on MAE and NSE alone to allow for greater focus on how models differ in
16 terms of error structure and uncertainty.

17 As a rough point of comparison for the statistical models developed in this research, we
18 also evaluated discharge estimates derived from a process-based hydrological model. The
19 model used in this application is the Noah Land Surface Model version 3.2 (Noah LSM; Ek
20 et. al, 2003; Chen et al., 1996). Noah LSM was implemented for offline simulations of the
21 Lake Tana basin at a gridded spatial resolution of 5km for the period 1979-2010 using a time
22 step of 30 minutes. Meteorological forcing was drawn from the Princeton 50-year reanalysis
23 dataset (Sheffield et al. 2006), downscaled to account for Ethiopia's steep terrain using
24 MicroMet elevation correction equations (Liston & Elder 2006). The Princeton reanalysis was
25 selected because it provides relatively high resolution meteorological fields, including all
26 variables required to run a water and energy balance LSM like Noah, for the period 1948-
27 present. While higher resolution and possibly higher quality datasets are available for recent
28 years, this longer dataset was utilized to compare the process-based model to statistical
29 models developed for a long historical period. Soil parameters for the Noah simulation were
30 drawn from the FAO global soil database, land use was defined according to the United States
31 Geological Survey (USGS) global 1km land cover product, and vegetation fraction was
32 derived from MODerate Imaging Spectroradiometer (MODIS) imagery. Land cover was

1 treated as a static parameter over the full length of the simulation, as spatially complete
2 estimates of historical land use were not available at the required resolution and specificity.

3 The highest performing model in each basin based on MAE was retained for more
4 detailed evaluation of model error structure, covariate influence, and uncertainty in climate
5 change sensitivity analysis. To generate a complete time-series of out-of-sample model
6 predictions for error analysis, the holdout cross validation procedure was repeated for the
7 highest performing standard-formulation and anomaly-formulation models for each basin, but
8 this time holding out a single year of observations in each iteration. The predictions from this
9 cross validation were used to evaluate the how model error structure might impact model
10 predictions used for water resource applications. The influence of different predictor variables
11 on model predictions was also assessed for the highest performing model in each basin after
12 being fit to the complete dataset. Each predictor variable was assessed using metrics for
13 covariate importance and influence that are unique to that model type, demonstrating how
14 models could be used to gain physical insights about data-scarce regions and the mechanisms
15 for generating these insights for each type of model. Partial dependence plots (Hastie et al.,
16 2009) were also generated for each covariate for the highest performing model in each basin
17 to provide insights about how covariate influence compared across different basins and model
18 types.

19 Finally, two evaluations were conducted to assess uncertainty in model projections of
20 streamflow under increasingly extreme climate conditions to better understand the
21 implications of using different model formulations for climate change impact studies. Model
22 projections of streamflow in different climate conditions are likely to be accompanied by
23 considerable uncertainty, particularly when climate conditions exceed those experienced
24 historically. To assess this uncertainty, the best performing model in each basin was used to
25 generate streamflow predictions for 1) changes in temperature from 0 to 5° C, 2) changes in
26 precipitation from -30 to +30%, 3) an increase in temperature to 5° C combined with a
27 decrease in precipitation to -30%, and 4) an increase in temperature to 5° C combined with an
28 increase in precipitation to +30%. For each of the four assessments, the models generated
29 predictions for the 45-year historic climate record adjusted for a given degree of climate
30 change using the delta-change method (Gleick, 1986), while holding agricultural land cover
31 constant at 60%. In this method, monthly temperature values are simply added to the
32 temperature change value, and monthly precipitation values are multiplied by the precipitation

1 change percentage. Model predictions for the altered climate record were then used to
2 calculate the average annual streamflow in each river. ~~These should not be interpreted as a~~
3 ~~prediction of actual climate change impacts, but rather a measurement of the sensitivity of~~
4 ~~streamflow in the basin to different climate conditions.~~ This process was repeated 100 times
5 for models fit on random bootstrap resamples of the historic dataset to generate uncertainty
6 bounds surrounding model predictions and evaluated how the uncertainty in these predictions
7 increased as climate conditions became more extreme. It is important to recognize that these
8 should not be interpreted as a prediction or assessment of actual climate change impacts, but
9 rather a measurement of the sensitivity of modeled streamflow in the basin to different
10 climate conditions. Since one of the key motivations for using rainfall-runoff models is to
11 understand how climate change may impact water resources, it is important to understand
12 how model formulation contributes to this sensitivity and uncertainty.

13 **3 Results**

14 **3.1 Model Accuracy and Error Structure**

15 Table 32 shows the out-of-sample cross validation errors for each model assessed in
16 each basin. The random forest model had the lowest mean absolute error for the standard-
17 formulation model in four of the five basins, with the M5 model performing best for the Koga
18 basin. These models outperformed the Noah LSM simulations in all basins assessed. The
19 Noah LSM errors are for a single period of analysis and thus don't present an exact corollary
20 to the cross validation performed for the empirical models. Nevertheless, the significant
21 increases in errors associated with the Noah LSM model demonstrates the difficulty
22 associated with the use of process-based models in the region, particularly when relying on
23 global datasets that may be unreliable at the spatial and temporal resolutions required for
24 physical modeling. Physical models developed for monthly streamflow prediction in other
25 basins within the Ethiopian highlands have reported NSE values ranging from 0.53 to 0.92
26 (van Griensven et al., 2012), compared to values ranging from 0.71 to 0.87 for the random
27 forest models developed here. If this measure alone was used for model evaluation, these
28 empirical models would generally be classified as having good performance based on the
29 guidelines suggested by Moraisi et al. (2007). However, the climatology model outperforms
30 the best standard formulation models in all basins except Megech, indicating that in the
31 majority of basins the errors from the fitted empirical models are higher than those that result
32 from simply using the long-term monthly average for each month's prediction. This is due to

1 the fact that seasonality accounts for such a large portion of the variability in monthly flow
2 values, and demonstrates how high NSE values can be quite easy to obtain in seasonal basins.

3 Evaluation of anomaly model errors indicates that the models using this formulation
4 achieve better predictive accuracy than those using the standard formulation, and are able to
5 outperform the climatology model based on both NSE and MAE in all basins. However, the
6 highest performing models in each basin varies more when the anomaly formulation is used,
7 with the GLM, GAM, random forest, and M5 models all minimizing MAE in different basins.
8 In all basins except Koga, the highest performing model significantly outperformed the
9 climatology model based on paired Wilcoxon rank-sum tests (Bonferroni-corrected p-value <
10 0.01).

11 Further exploration of model residuals indicates another important advantage of using
12 the anomaly model formulation. In the standard model formulation, model residuals appear to
13 be non-random. Example autocorrelation plots are shown for the Gilgel Abbay and Ribb
14 Rivers in Fig. 2, and demonstrate that a positive autocorrelation exists at the 12 month time
15 lag. For brevity, only plots for two rivers are shown, although this autocorrelation existed in
16 the standard-formulation models for all basins except Megech (Table 43). This
17 autocorrelation occurs because the standard-formulation models consistently underestimate
18 wet-season streamflow while overestimating dry-season flows, as is apparent in hydrographs
19 of observed and predicted streamflow (Fig. 3). Because wet-season flows contribute such a
20 large portion of the total annual flow volume, this results in regular underestimation of
21 aggregate values such as mean annual flow (Table 43). This autocorrelation is reduced in the
22 anomaly-formulation models, meaning that they are better able to capture the peak flow
23 volumes experienced in the wet season and do not underestimate mean annual flow to the
24 same degree that the standard formulation models do.

25 **3.2 Model Structure and Covariate Influence**

26 Evaluating the relationship between predictor covariates and streamflow response can
27 lend insight into the physical processes underlying runoff generation in each basin. There are
28 two components of this relationship that can be evaluated: how much each covariate
29 contributes to model accuracy (covariate importance), and the direction and nature of the
30 relationship between covariate values and model response (covariate influence). In many
31 machine-learning models, complete description of the all of the mathematical relationships

1 within the model (for instance, through description of each tree comprising a random forest
2 model) is infeasible, requiring the use of other mechanisms for understanding covariate
3 importance and influence. However, because each model type is structured in a different way,
4 these mechanisms differ. This section first describes the mechanisms available for obtaining
5 insights about covariate influence in each of the highest performing models. To provide a
6 mechanism for comparing results across different basins, each basin model is then assessed
7 using the general approach of partial dependence plots.

8 In the Gilgel Abbay and Koga basins, the highest performing model was a simple
9 linear regression model. These models can be evaluated by reviewing model coefficients and
10 associated p-values, as shown in Table 54. In a standard linear regression, model coefficients
11 can be interpreted as the mean change in the response variable that results from a unit change
12 in that covariate when all others are held constant. These coefficients are for streamflow
13 anomalies rather than raw values, making their immediate interpretation less intuitive. For
14 instance, in the Gilgel Abbay model an increase of one standard deviation in precipitation
15 results in an increase of 0.22 standard deviations in flow. The associated p-value for each
16 coefficient evaluates a null hypothesis that the true coefficient value is equal to zero given the
17 other covariates in the model, and thus has no influence on the response variable.

18 Evaluating model structure based on regression coefficients is appealing due to their
19 simplicity and familiarity. However, it is important to keep in mind that the above
20 interpretations rely on specific assumptions regarding model error distributions. Examination
21 of fitted model residuals from both basins indicate that errors are autocorrelated in the Koga
22 basin and not normally distributed due to the presence of outliers in both basins. Non-
23 normality and autocorrelation both impact the t statistics and f statistics used to test for the
24 significance of model coefficients, and thus the p-values for these models are likely biased
25 (Montgomery et al., 2012).

26 Interpretation of variable influence in GAMs is based on the estimated degrees of
27 freedom (EDF) a covariate's smoothing function $s(X_i)$ uses within a model (Hastie and
28 Tibushini, 1986). An EDF value of one or below indicates a linear function relating the
29 response variable to that covariate, while values greater than one represent a non-linear
30 smoothing function. An EDF value of zero indicates that the covariate smoothing function is
31 penalized to zero (meaning it has no influence on model predictions). In the model for the
32 Megech River, the terms for lagged temperature at one and two months, as well as

1 precipitation lagged at two months were all smoothed to zero. Of the remaining covariates,
2 lagged precipitation has a linear impacts on model response, while precipitation, temperature
3 and land cover have non-linear impacts. Smoothing functions can be plotted to gain more
4 insight about these relationships (Fig. 4). The functions for precipitation anomaly, lagged (one
5 month) precipitation anomaly, and agricultural land cover show a positive relationships with
6 streamflow, while the function for temperature anomaly predicts low streamflow at both high
7 and low anomalies.

8 P-values test the null hypothesis that a covariate's smoothing function is equal to zero,
9 but rest on the assumption that model residuals are homoscedastic and independent (Wood,
10 2012). Similar to the linear models, residuals in the Megech GAM model appear to be both
11 autocorrelated and heteroscedastic, meaning that a formal statistical interpretation of this
12 value may be inappropriate and that confidence bounds around smoothing functions might be
13 misleading.

14 The M5 cubist model fit for the Gumara basin is an ensemble of 100 small M5
15 regression trees. In each tree, the model splits observations based on logical rules related to
16 one or more covariates and fits a linear regression model to each set of observations. The final
17 model prediction is the average across all of the individual trees. Using this sort of ensemble
18 approach can reduce model variance and improve accuracy if the individual trees are
19 unbiased, uncorrelated predictors (Breiman 1996). This can be useful in avoiding models that
20 are overfit to the data, but can reduce model interpretability since direct visualization of
21 model structure becomes impractical as the number of trees increases. However, the
22 frequency with which individual covariates are used as splitting points within trees and as
23 regression coefficients can provide some insights about covariate importance (Table 54; note
24 that because multiple covariates can be used for rules and linear models, these don't
25 necessarily add to 100%). Model rules were largely based on land cover, with some rules
26 based on precipitation. These two covariates were also used most frequently in linear
27 regressions at model nodes, followed by temperature (current and 1-month lag) and 1-month
28 lagged precipitation. Notably, climate data from 2 months lagged ~~were~~ not used at all.
29 While this can be useful in identifying which covariates have the largest impact on model
30 predictions, it doesn't provide any information regarding the nature or direction of that
31 influence.

1 Similarly, the random forest model developed for the Ribb basin is an ensemble of
2 regression trees in which the final model prediction is the average of the predictions from
3 each individual tree. However, random forests use standard regression trees that do not
4 incorporate linear regression models at terminal nodes. Variable importance within the final
5 model is measured by recording the increase in out-of-sample MSE that results when a
6 covariate is randomly permuted for each tree in the ensemble. This increase in error is then
7 averaged across all trees in the ensemble. In our model, the largest increases in error resulted
8 from permutation of land cover and temperature, followed by 2-month lagged temperature
9 and precipitation. Covariate influence can be evaluated through the use of partial dependence
10 plots, which measure the change in model predictions that result from changing the value of
11 one parameter while leaving all other covariates constant (Hastie et al., 2009). Partial
12 dependence plots indicate that model predictions of streamflow are higher when the percent of
13 agricultural land cover is greater than approximately 75%, when temperatures anomalies are
14 low, and when precipitation anomalies are high. However, it appears that the plot for lagged
15 temperature might be sensitive to outliers at high temperature anomalies as evidenced by the
16 large increase that occurs above an anomaly of +2, in a region where very few data points are
17 present.

18 Many of the measures used to evaluate covariate importance and influence are model
19 specific, making inter-basin and inter-model comparisons difficult. However, the partial
20 dependence plots used in the randomForest R package can be developed for any model and
21 provide a mechanism for comparing the influence that covariates have in the different models
22 and basins (Shortridge et al., 2015). Partial dependence plots were generated for each basin's
23 best performing model and results are shown for climatic variables in Fig. 6. As expected,
24 models generally respond positively to increases in precipitation and negatively to increases
25 in temperature, with the greatest influence in the current month and decreasing influence at
26 one and two months prior. The influence of the current month's precipitation is linear in three
27 of the five basins; while this is constrained to be the case in the Gilgel Abbay and Koga
28 basins due to the use of a linear model, the linear response in Gumara is not required from the
29 M5 model structure. Interestingly, both Megech and Ribb demonstrate a linear response to
30 negative precipitation anomalies, but little response to positive anomalies. Streamflow
31 response to temperature is strongest in the Gumara basin; interestingly, this is the basin with
32 the smallest response to precipitation.

1 The partial dependence plots for the percentage of the basin classified as agricultural
2 land cover indicates a positive relationship between agricultural land cover and streamflow in
3 all basins except for the Gilgel Abbay (Fig. 7). This would be expected if deforestation had
4 contributed to a decrease in evapotranspiration in the contributing watersheds. The exact
5 nature of this response differs across the different rivers, with the relatively minor responses
6 in Koga and Ribb, and much stronger responses in the Gumara and Megech basins. However,
7 this plot also demonstrates some of the limitations associated with different model structures.
8 The plot for Gumara is highly erratic, indicating that the M5 model might be overfit to the
9 training dataset, despite the use of model averaging to reduce model variance. Additionally,
10 the GAM used in the Megech basin was only trained on agricultural land cover values up to
11 77%; while this model may be accurately representing the impact of land cover changes
12 within this range, extrapolating this relationship to higher values leads to predictions that may
13 not be physically realistic.

14 **3.3 Climate Change Sensitivity and Uncertainty Assessment**

15 Fig. 8 shows the results of the climate change sensitivity analysis for total flow from all
16 five tributaries, with dashed lines representing 95% confidence intervals obtained through 100
17 bootstrapped resamples of the data set. As would be expected, increasing temperature
18 independently of precipitation results in decreasing total flows while increasing precipitation
19 results in higher flows. However, the uncertainty surrounding temperature sensitivity
20 increases at higher changes in temperature, while the uncertainty surrounding precipitation
21 sensitivity remains relatively constant, even at extreme changes in annual precipitation. The
22 bottom panels of the figure show the sensitivity of total inflows to concurrent changes in
23 temperature and precipitation. Unsurprisingly, decreasing precipitation combined with higher
24 temperatures results in greater decreases in total flow than when temperature and precipitation
25 are varied independently. However, even if temperature increases are combined with higher
26 precipitation, total flows decline in the majority of bootstrap resamples.

27 The uncertainty surrounding temperature sensitivity is a key limitation to using data-
28 driven approaches for climate impact assessment. To better understand which models and
29 basins are contributing to this uncertainty, Fig. 9 shows how the coefficient of variation (the
30 standard deviation of predictions from all bootstrap samples divided by the mean of these
31 predictions) varies as a function of temperature change in each basin. From this figure, it is
32 apparent that the Megech model is by far the largest contributor to model uncertainty;

1 however, it is not clear whether this contribution is due to model structure (the GAM model
2 used for the Megech River) or characteristics associated with the basin itself. To investigate
3 how different model structures contributed to this uncertainty, the bootstrap resampling
4 procedure was used to assess uncertainty in streamflow predictions in the Gumara River from
5 all model types. This basin was chosen because all six models were able to outperform the
6 climatology model, and thus could be considered good choices for model selection based on
7 predictive accuracy alone. The results indicate that the increase in uncertainty is highest, and
8 increases non-linearly, in the GLM, GAM, and MARS models. Uncertainty increases more
9 slowly in the ANN and M5 models, and no noticeable increase in uncertainty is apparent in
10 the random forest model.

11 **4 Discussion**

12 The objective of this study was not to identify the “best” approach for empirical
13 rainfall-runoff modeling, as this is likely to be highly specific to the basin and problem to
14 which a model is applied. However, we hope that the comparison conducted here can
15 highlight some of the strengths and limitations of different approaches, as well as demonstrate
16 some important issues that should be kept in mind for model comparisons in the future. One
17 important finding was the limitation with using NSE as an error metric. Our results confirm
18 previous studies that found that even uninformative models able to capture basic seasonality
19 are able to achieve high NSE values (Legates and McCabe, 1999; Schaefli and Gupta, 2007),
20 and provide further evidence indicating that high NSE values should be considered a
21 necessary but not sufficient requirement for model usage in planning situations. For instance,
22 the simple climatology model used for comparison purposes here is able to achieve high NSE
23 values, but would be unsuitable for planning since it does not account for any interannual
24 variability nor the possibility for non-stationary conditions caused by changing climate and
25 land cover. In particular, understanding error structure can be valuable in evaluating whether
26 model biases might undermine the model’s suitability for management activities. In our
27 example, the autocorrelation present in the standard-formulation models meant that these
28 models were consistently underestimating wet-season flows, resulting in low estimates of the
29 total annual flow in the rivers. Since multiple reservoirs are planned for construction on these
30 rivers to support irrigation activities, this bias could lead to poor estimates of how much water
31 is available for agricultural use in the short term (ie., seasonal forecasting) and long-term (due
32 to climate change). Interestingly, difficulties in accurately capturing high flows has been

1 observed in physical hydrologic models for Ethiopia (e.g., Setegne et al., 2011; Mekonnen et
2 al., 2009) and more generally (e.g., Wilby, 2005). The implications of this limitation should
3 be carefully evaluated before using models for water resource planning or (more importantly)
4 flood risk evaluation.

5 Depending on the model type used, different mechanisms are available to evaluate
6 covariate importance and influence within the model. This evaluation can be useful in ~~both~~
7 confirming that the model is replicating physically realistic relationships between input and
8 output variables, ~~and in characterizing watershed behavior in a manner that could shed light~~
9 ~~on underlying physical processes.~~ While the relationships identified in this evaluation are
10 fairly straightforward (for example, increasing runoff with higher precipitation and lower
11 temperatures), these simple relationships are still important in highlighting the mechanisms by
12 which the models make predictions so that they are not “black boxes.” For instance, Han et al.
13 (2007) explore how ANN flood forecasting models responds to a double-unit input of rain,
14 finding that some formulations respond in a hydrologically meaningful way to increased
15 rainfall intensity, while others do not. Similarly, Galelli and Castelletti (2013a) describe how
16 input variable importance can be used to highlight differences in hydrologic processes
17 between an urbanized and forested watershed. The easy manner in which covariate
18 relationships within the GAM and random forest models can be visualized using a single
19 command within their respective R packages is a strong advantage to these approaches
20 compared to methods such as M5 model trees and artificial neural networks. Of course, partial
21 dependence plots can be developed for any model type (as was done in this research), but
22 code must be written by the user and thus requires a higher degree of effort than is necessary
23 for in-package functions. A downside to most machine-learning models is that they do not
24 support the statistical formalism in assessing variable importance that is possible when linear
25 models and GAMs are used. However, this formalism often rests on assumptions regarding
26 model residuals that are unlikely to be met in many hydrologic models (Sorooshian and
27 Dracup, 1980).

28 Within the Lake Tana basin, evaluation of covariate influence indicates that each
29 basin’s model is performing in a physically realistic manner, with a runoff increasing with
30 higher precipitation levels and decreasing with higher temperatures. The influence of
31 precipitation and temperature is greatest in the current month, and progressively declines to a
32 very small influence after two months. This suggests that long-term (multi-month) storage

1 does not significantly contribute to variability in flow volumes. One interesting finding is the
2 non-linear relationship between concurrent month precipitation and runoff that exists in the
3 Megech and Ribb basins, which suggests that above a certain point increasing rainfall does
4 not result in a commiserate increase in streamflow. Other studies have noted the dampening
5 effect that wetlands and floodplains have had on river flows in the region (Dessie et al., 2014;
6 Gebrehiwot et al., 2010); this phenomenon could explain the non-linear relationship identified
7 in this work. The clearly negative relationship between temperature and runoff demonstrates
8 the degree to which upstream evapotranspiration impacts streamflow and suggests that
9 evapotranspiration is largely energy-limited, rather than water-limited. Increasing agricultural
10 land-use appears to be associated with higher runoff in all rivers except for Gilgel Abbay
11 (where no clear relationship between land cover and runoff was observed), and suggests that
12 agricultural expansion at the expense of forest cover has reduced the evaporative component
13 of the water balance in these basins. Finally, the relative performance of different model
14 formulations themselves can also be informative. For instance, the improved performance of
15 the anomaly-formulation models indicates that the relationship between precipitation and
16 runoff varies throughout the year and could point towards differences in runoff-generating
17 mechanisms in the wet and dry seasons that have been observed in other case studies (Wilby,
18 2005).

19 One limitation with data-driven approaches for streamflow prediction is that the
20 relationships they model can only generate reliable predictions for conditions that are
21 comparable to those experienced historically. Using these models to generate predictions for
22 conditions that exceed historic variability is likely to introduce considerable uncertainty into
23 their projections. Our results indicate that uncertainty in projections of streamflow under
24 changing precipitation is relatively constant, whereas uncertainty increases markedly in
25 projections of streamflow under increasing temperature. This result is not surprising when one
26 considers the basin's climate, which is characterized by highly variable rainfall but fairly
27 consistent temperatures (Table ~~65~~). A temperature increase of 3° C equates to almost two
28 standard deviations beyond historic variability, whereas a change in precipitation of 30% is
29 well within the range of conditions experienced historically. One would expect that in other
30 climates (for example, temperate watersheds with only minor changes in rainfall throughout
31 the year), this relationship could be reversed. Despite the uncertainty that exists in projections
32 of streamflow under changing temperature, total annual flow appears to be quite sensitive to
33 increasing temperatures. In fact, the decreases in streamflow due to increasing temperature

1 appears likely to be more than enough to counteract any increases in streamflow resulting
2 from higher precipitation that is projected for the region in some global circulation models
3 (GCMs). This is consistent with the work of Setegne et al. (2011), who used projections from
4 multiple GCMs as input for a SWAT model developed for the region and found that
5 streamflow decreased in the majority of emissions scenarios and models, even when
6 precipitation increased. Unfortunately, this suggests that any hopes for a “windfall” of
7 additional water to support agriculture and hydropower in the region under climate change
8 may be unfounded.

9 Repeating the climate change sensitivity experiment with multiple models fit to the
10 Gumara watershed indicated that the MARS, GAM, and linear models all result in the largest
11 increase in uncertainty at high temperatures. This indicates that when models are fit to slightly
12 different bootstrap resamples of the historic dataset, the projected changes in streamflow at
13 high temperature changes can be highly erratic. This is likely due to the fact that extrapolating
14 the relationships that are observed between historic temperature and streamflow to higher
15 temperatures can lead to very large changes in streamflow. Fitting the models to bootstrap
16 resamples of the data results in minor changes to these relationships that can result in widely
17 varying projections when the models are used to predict streamflow at higher temperatures,
18 particularly when these relationships are nonlinear (as in the GAM). At the other end of the
19 spectrum, the random forest model exhibits almost no increase in uncertainty at high
20 temperatures, meaning that projections of streamflow at high temperatures are consistent
21 across the bootstrap resamples. This is likely the result of the random forest model structure.
22 The predicted value for each of a regression tree’s terminal nodes is the average of all
23 observations that meet the conditions described for that node. Thus, the model will not predict
24 values beyond those experienced historically, even if covariate values exceed those contained
25 within the historic dataset. Thus, this model is likely to underestimate the change in
26 streamflow that results from increasing temperatures.

27 **5 Conclusions**

28 In this work, we compared multiple methods for data-driven rainfall-runoff modeling
29 in their ability to simulate streamflow in five highly-seasonal watersheds in the Ethiopian
30 highlands. Despite the popularity of ANNs in research on streamflow prediction to date,
31 ANNs were not found to be the most accurate model in any of the five basins evaluated. Other
32 methods, in particular GAMs and random forests, are able to capture non-linear relationships

1 effectively and lend themselves to simpler visualization of model structure and covariate
2 influence, making it easier to gain insights on physical watershed functions and confirm that
3 the model is operating in a physically realistic manner. However, it is important to carefully
4 evaluate model structure and residuals, as these can contribute to biased estimates of water
5 availability and uncertainty in estimating sensitivity to potential future changes in climate. In
6 particular, autocorrelation in model residuals can result in underestimation of aggregate
7 metrics such as annual flow volumes, even in models with high NSE performance.
8 Uncertainty in GAM projections was found to rapidly increase at high temperatures, whereas
9 random forest projections may be underestimating the impact of high temperatures on river
10 flows. Thorough consideration of this uncertainty and bias is important any time that models
11 are used for water planning and management, but especially crucial when using such models
12 to generate insights about future streamflow levels. By considering multiple model
13 formulations and carefully assessing their predictive accuracy, error structure and
14 uncertainties, these methods can provide an empirical assessment of watershed behavior and
15 generate useful insights for water management and planning. This makes them a valuable
16 complement to physical models, particularly in data-scarce regions with little data available
17 for model parameterization, and warrants additional research into their development and
18 application.

19 **Acknowledgements**

20 We would like to gratefully acknowledge the Ethiopian Ministry of Water and Energy, the
21 Tana Sub Basin Organization, and the International Water Management Institute for making
22 available the data used to perform this analysis. All data for this paper [are](#) properly cited and
23 referred to in the reference list. The source code for the models developed in this study is
24 available from the authors upon request. Empirical modeling work was supported by a
25 National Defense Science and Engineering Graduate Fellowship and by National Science
26 Foundation Grant 1069213 (IGERT). Noah LSM simulations presented here were performed
27 under NASA Applied Sciences Program grant NNX09AT61G. This research was conducted
28 while Dr. Guikema was affiliated with the Department of Geography and Environmental
29 Engineering at Johns Hopkins University. This support is gratefully acknowledged. Any
30 opinions, findings, and conclusions or recommendations expressed in this material are those
31 of the authors and do not necessarily reflect the views of the funding sources.

32

33

1 **References**

- 2 Abraham, R. J. and See, L. M.: Neural network modelling of non-linear hydrological
3 relationships, *Hydrol. Earth Syst. Sci.*, 11(5), 1563–1579, doi:10.5194/hess-11-1563-2007,
4 2007.
- 5 Achenef, H., Tilahun, A. and Molla, B.: Tana Sub Basin Initial Scenarios and Indicators
6 Development Report, Tana Sub Basin Organization, Bahir Dar, Ethiopia., 2013.
- 7 Alemayehu, T., McCartney, M. and Kebede, S.: The water resource implications of planned
8 development in the Lake Tana catchment, Ethiopia, *Ecohydrology & Hydrobiology*, 10(2-4),
9 211–221, doi:10.2478/v10104-011-0023-6, 2010.
- 10 [Antar, M. A., Ellassiouti, I. and Allam, M. N.: Rainfall-runoff modelling using artificial neural
11 networks technique: a Blue Nile catchment case study, *Hydrol. Process.*, 20\(5\), 1201–1216,
12 doi:10.1002/hyp.5932, 2006.](#)
- 13 [Aqil, M., Kita, I., Yano, A. and Nishiyama, S.: Neural Networks for Real Time Catchment
14 Flow Modeling and Prediction, *Water Resour. Manag.*, 21\(10\), 1781–1796, doi:
15 10.1007/s11269-006-9127-y, 2007.](#)
- 16 Asefa, T., Kemblowski, M., McKee, M. and Khalil, A.: Multi-time scale stream flow
17 predictions: The support vector machines approach, *J. Hydrol.*, 318(1-4), 7–16,
18 doi:10.1016/j.jhydrol.2005.06.001, 2006.
- 19 Beven, K. J.: *Rainfall-Runoff Modelling: The Primer*, John Wiley & Sons., 2011.
- 20 Breiman, L.: Bagging predictors, *Mach Learn*, 24(2), 123–140, doi:10.1007/BF00058655,
21 1996.
- 22 Breiman, L.: Random forests, *Mach Learn*, 45(1), 5–32, 2001.
- 23 Chen, F., Mitchell, K., Schaake, J., Xue, Y., Pan, H.-L., Koren, V., Duan, Q. Y., Ek, M. and
24 Betts, A.: Modeling of land surface evaporation by four schemes and comparison with FIFE
25 observations, *J. Geophys. Res.*, 101(D3), 7251–7268, doi:10.1029/95JD02165, 1996.
- 26 [Chibanga, R., Berlamont, J. and Vandewalle, J.: Modelling and forecasting of hydrological
27 variables using artificial neural networks: the Kafue River sub-basin, *Hydrolog. Sci. J.*, 48\(3\),
28 363–379, doi:10.1623/hysj.48.3.363.45282, 2003.](#)
- 29 [Criss, R. E. and Winston, W. E.: Do Nash values have value? Discussion and alternate](#)

1 [proposals, Hydrol. Process., 22\(14\), 2723–2725, doi:10.1002/hyp.7072, 2008.](#)

2 [De Vos, N. J. and Rientjes, T. H. M.: Multiobjective training of artificial neural networks for](#)
3 [rainfall-runoff modeling, Water Resour. Res., 44\(8\), W08434, doi:10.1029/2007WR006734,](#)
4 [2008.](#)

5 Dessie, M., Verhoest, N. E. C., Admasu, T., Pauwels, V. R. N., Poesen, J., Adgo, E., Deckers,
6 J. and Nyssen, J.: Effects of the floodplain on river discharge into Lake Tana (Ethiopia), J.
7 Hydrol., 519, 699–710, doi:10.1016/j.jhydrol.2014.08.007, 2014.

8 Ek, M. B., Mitchell, K. E., Lin, Y., Rogers, E., Grunmann, P., Koren, V., Gayno, G. and
9 Tarpley, J. D.: Implementation of Noah land surface model advances in the National Centers
10 for Environmental Prediction operational mesoscale Eta model, J. Geophys. Res., 108(D22),
11 8851, doi:10.1029/2002JD003296, 2003.

12 Elshorbagy, A., Corzo, G., Srinivasulu, S. and Solomatine, D. P.: Experimental investigation
13 of the predictive capabilities of data driven modeling techniques in hydrology - Part 1:
14 Concepts and methodology, Hydrol. Earth Syst. Sci., 14(10), 1931–1941, doi:10.5194/hess-
15 14-1931-2010, 2010a.

16 Elshorbagy, A., Corzo, G., Srinivasulu, S. and Solomatine, D. P.: Experimental investigation
17 of the predictive capabilities of data driven modeling techniques in hydrology - Part 2:
18 Application, Hydrol. Earth Syst. Sci., 14(10), 1943–1961, doi:10.5194/hess-14-1943-2010,
19 2010b.

20 Friedman, J. H.: Multivariate adaptive regression splines, The annals of statistics, 1–67, 1991.

21 Galelli, S. and Castelletti, A.: Assessing the predictive capability of randomized tree-based
22 ensembles in streamflow modelling, Hydrol. Earth Syst. Sci., 17(7), 2669–2684,
23 doi:10.5194/hess-17-2669-2013, 2013.

24 [Galelli, S. and Castelletti, A.: Tree-based iterative input variable selection for hydrological](#)
25 [modeling, Water Resour. Res., 49\(7\), 4295–4310, doi:10.1002/wrcr.20339, 2013.](#)

26 Garede, N. M. and Minale, A. S.: Land Use/Cover Dynamics in Ribb Watershed, North
27 Western Ethiopia, Journal of Natural Sciences Research, 4(16), 9–16, 2014.

28 Gaume, E. and Gosset, R.: Over-parameterisation, a major obstacle to the use of artificial
29 neural networks in hydrology?, Hydrol. Earth Syst. Sci. Discussions, 7(5), 693–706, 2003.

30 Gebrehiwot, S. G., Taye, A. and Bishop, K.: Forest Cover and Stream Flow in a Headwater of

1 the Blue Nile: Complementing Observational Data Analysis with Community Perception,
2 *Ambio*, 39(4), 284–294, doi:10.1007/s13280-010-0047-y, 2010.

3 Han, D., Kwong, T. and Li, S.: Uncertainties in real-time flood forecasting with neural
4 networks, *Hydrol. Process.*, 21(2), 223–228, doi:10.1002/hyp.6184, 2007.

5 Harris, I., Jones, P. d., Osborn, T. j. and Lister, D. h.: Updated high-resolution grids of
6 monthly climatic observations – the CRU TS3.10 Dataset, *Int. J. Climatol.*, 34(3), 623–642,
7 doi:10.1002/joc.3711, 2014.

8 Hastie, T. and Tibshirani, R.: *Generalized Additive Models*, *Statistical Science*, 1(3), 297–
9 310, 1986.

10 Hastie, T. and Tibshirani, R.: *Generalized additive models*. Chapman, Hall, London, 1990.

11 Hastie, T., Tibshirani, R. and Friedman, J.: *The Elements of Statistical Learning: Data
12 Mining, Inference and Prediction*, Second Ed., Springer, New York., 2009.

13 Iorgulescu, I. and Beven, K. J.: Nonparametric direct mapping of rainfall-runoff relationships:
14 An alternative approach to data analysis and modeling?, *Water Resour. Res.*, 40(8), W08403,
15 doi:10.1029/2004WR003094, 2004.

16 Jain, A., Sudheer, K. P. and Srinivasulu, S.: Identification of physical processes inherent in
17 artificial neural network rainfall runoff models, *Hydrol. Process.*, 18(3), 571–581,
18 doi:10.1002/hyp.5502, 2004.

19 Kuhn, M.: caret: Classification and regression training, Available from: [http://CRAN.R-](http://CRAN.R-project.org/package=caret)
20 [project.org/package=caret](http://CRAN.R-project.org/package=caret), 2015.

21 Kuhn, M., Weston, S., Keefer, C. and Coulter, N.: Cubist: Rule- and instance-based
22 regression modeling, Available from: <http://CRAN.R-project.org/package=Cubist>, 2014.

23 Legates, D. R. and McCabe Jr, G. J.: Evaluating the use of“ goodness-of-fit” measures in
24 hydrologic and hydroclimatic model validation, *Water Resour. Res.*, 35(1), 233–241, 1999.

25 Liaw, A. and Wiener, M.: Classification and regression by randomForest, *R News*, 2(3), 18–
26 22, 2002.

27 Lin, J.-Y., Cheng, C.-T. and Chau, K.-W.: Using support vector machines for long-term
28 discharge prediction, *Hydrological Sciences Journal*, 51(4), 599–612,
29 doi:10.1623/hysj.51.4.599, 2006.

1 Liston, G. E. and Elder, K.: A Meteorological Distribution System for High-Resolution
2 Terrestrial Modeling (MicroMet), *J. Hydrometeor.*, 7(2), 217–234, doi:10.1175/JHM486.1,
3 2006.

4 [Machado, F., Mine, M., Kaviski, E. and Fill, H.: Monthly rainfall–runoff modelling using](#)
5 [artificial neural networks, *Hydrolog. Sci. J.*, 56\(3\), 349–361,](#)
6 [doi:10.1080/02626667.2011.559949, 2011.](#)

7 [Maier, H. R., Jain, A., Dandy, G. C. and Sudheer, K. P.: Methods used for the development of](#)
8 [neural networks for the prediction of water resource variables in river systems: Current status](#)
9 [and future directions, *Environ. Modell. Softw.*, 25\(8\), 891–909,](#)
10 [doi:10.1016/j.envsoft.2010.02.003, 2010.](#)

11 [Mathevet, T., Michel, C., Andreassian, V. and Perrin, C.: A bounded version of the Nash-](#)
12 [sutcliffe criterion for better model assessment on large sets of basins, in IAHS-AISH](#)
13 [publication, pp. 211–219, International Association of Hydrological Sciences. \[online\]](#)
14 [Available from: <http://cat.inist.fr/?aModele=afficheN&cpsidt=18790113> \(Accessed 10](#)
15 [February 2016\), 2006.](#)

16 Mekonnen, M. A., Wörman, A., Dargahi, B. and Gebeyehu, A.: Hydrological modelling of
17 Ethiopian catchments using limited data, *Hydrol. Process.*, 23(23), 3401–3408,
18 doi:10.1002/hyp.7470, 2009.

19 Milborrow, S.: earth: Multivariate Adaptive Regression Splines, Available from:
20 <http://CRAN.R-project.org/package=earth>, 2015.

21 Montgomery, D. C., Peck, E. A. and Vining, G. G.: Introduction to Linear Regression
22 Analysis, John Wiley & Sons., 2012.

23 Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D. and Veith, T.
24 L.: Model evaluation guidelines for systematic quantification of accuracy in watershed
25 simulations, *Trans. Asabe*, 50(3), 885–900, 2007.

26 [Pushpalatha, R., Perrin, C., Moine, N. L. and Andréassian, V.: A review of efficiency criteria](#)
27 [suitable for evaluating low-flow simulations, *J. Hydrol.*, 420–421, 171–182,](#)
28 [doi:10.1016/j.jhydrol.2011.11.055, 2012.](#)

29 Quinlan, J. R.: Learning with Continuous Classes, in Proceedings of the 5th Australian Joint
30 Conference on Artificial Intelligence, World Scientific, Singapore., 1992.

1 R Development Core Team: R: A language and environment for statistical computing., R
2 Foundation for Statistical Computing, Vienna, Austria. Available from: [http://www.R-](http://www.R-project.org)
3 [project.org](http://www.R-project.org), 2014.

4 Rientjes, T. H. M., Haile, A. T., Kebede, E., Mannaerts, C. M. M., Habib, E. and Steenhuis,
5 T. S.: Changes in land cover, rainfall and stream flow in Upper Gilgel Abbay catchment, Blue
6 Nile basin – Ethiopia, *Hydrol. Earth Syst. Sci.*, 15(6), 1979–1989, doi:10.5194/hess-15-1979-
7 2011, 2011.

8 Ripley, B. D.: *Pattern Recognition and Neural Networks*, Cambridge University Press., 1996.

9 Schaeffli, B. and Gupta, H. V.: Do Nash values have value?, *Hydrol. Process.*, 21(15), 2075–
10 2080, doi:10.1002/hyp.6825, 2007.

11 See, L., Solomatine, D., Abrahart, R., and Toth, E.: Hydroinformatics: computational
12 intelligence and technological developments in water science applications—Editorial,
13 *Hydrological Sciences Journal*, 52(3), 391–396, doi:10.1623/hysj.52.3.391, 2007.

14 Setegn, S. G., Srinivasan, R., Melesse, A. M. and Dargahi, B.: SWAT model application and
15 prediction uncertainty analysis in the Lake Tana Basin, Ethiopia, *Hydrol. Process.*,
16 doi:10.1002/hyp.7457, 2009.

17 Setegn, S. G., Rayner, D., Melesse, A. M., Dargahi, B. and Srinivasan, R.: Impact of climate
18 change on the hydroclimatology of Lake Tana Basin, Ethiopia, *Water Resour. Res.*, 47(4),
19 doi:10.1029/2010WR009248, 2011.

20 Sheffield, J., Goteti, G. and Wood, E. F.: Development of a 50-Year High-Resolution Global
21 Dataset of Meteorological Forcings for Land Surface Modeling, *J. Climate*, 19(13), 3088–
22 3111, doi:10.1175/JCLI3790.1, 2006.

23 Shortridge, J. E., Falconi, S. M., Zaitchik, B. F. and Guikema, S. D.: Climate, agriculture, and
24 hunger: statistical prediction of undernourishment using nonlinear regression and data-mining
25 techniques, *Journal of Applied Statistics* (ahead of press),
26 doi:10.1080/02664763.2015.1032216, 2015.

27 Solomatine, D. P. and Ostfeld, A.: Data-driven modelling: some past experiences and new
28 approaches, *Journal of Hydroinformatics*, 10(1), 3, doi:10.2166/hydro.2008.015, 2008.

29 Sorooshian, S. and Dracup, J. A.: Stochastic parameter estimation procedures for hydrologic
30 rainfall-runoff models: Correlated and heteroscedastic error cases, *Water Resour. Res.*, 16(2),

1 430–442, doi:10.1029/WR016i002p00430, 1980.

2 Steenhuis, T. S., Collick, A. S., Easton, Z. M., Leggesse, E. S., Bayabil, H. K., White, E. D.,
3 Awulachew, S. B., Adgo, E. and Ahmed, A. A.: Predicting discharge and sediment for the
4 Abay (Blue Nile) with a simple model, *Hydrol. Process.*, doi:10.1002/hyp.7513, 2009.

5 Sudheer, K. P. and Jain, A.: Explaining the internal behaviour of artificial neural network
6 river flow models, *Hydrol. Process.*, 18(4), 833–844, doi:10.1002/hyp.5517, 2004.

7 Van Griensven, A., Ndomba, P., Yalew, S. and Kilonzo, F.: Critical review of SWAT
8 applications in the upper Nile basin countries, *Hydrol. Earth Syst. Sci.*, 16(9), 3371–3381,
9 doi:10.5194/hess-16-3371-2012, 2012.

10 Venables, W. N. and Ripley, B. D.: *Modern Applied Statistics with S-PLUS*, Springer
11 Science & Business Media., 2013.

12 Wilby, R. L.: Uncertainty in water resource model parameters used for climate change impact
13 assessment, *Hydrol. Process.*, 19(16), 3201–3219, doi:10.1002/hyp.5819, 2005.

14 Wood, S.: *Generalized Additive Models: An Introduction with R*, CRC Press., 2006.

15 Wood, S. N.: Fast stable restricted maximum likelihood and marginal likelihood estimation of
16 semiparametric generalized linear models, *Journal of the Royal Statistical Society: Series B*
17 *(Statistical Methodology)*, 73(1), 3–36, doi:10.1111/j.1467-9868.2010.00749.x, 2011.

18 Wood, S. N.: On p-values for smooth components of an extended generalized additive model,
19 *Biometrika*, 100(1), 221–228 doi:10.1093/biomet/ass048, 2012.

20

21

1 Table 1. Study basin characteristics over the evaluation period of 1961 to 2004.

Basin	Drainage area above gauge (km ²)	Average annual streamflow at gauge (MCM)	Standard deviation of annual streamflow (MCM)	<u>Coefficient of variation of annual streamflow</u>	Average temp (°C)	Average monthly rainfall [mm]	
						May-Oct	Nov-Apr
Gilgel Abbay	2664	1883	217	<u>0.12</u>	15.7	206	39.3
Gumara	385	236	71	<u>0.30</u>	17.7	186	29
Koga	200	114	31	<u>0.27</u>	15.7	206	39.3
Megech	424	172	66	<u>0.31</u>	20.6	234	41.4
Ribb	677	210	83	<u>0.36</u>	18.2	263	45.8

2

1 Table 2. Model parameters evaluated through cross validation.

Model type	R package	Parameters defined in model formulation	Parameters selected through cross validation
GLM	stats	<u>family = Gaussian</u>	NA
GAM	mgcv	<u>family = Gaussian</u> <u>method = generalized cross validation</u> <u>variable selection = true</u> <u>basis dimension k = 3</u> <u>epsilon = 10⁻⁷</u> <u>maxit = 200</u>	
MARS	earth	<u>nk = 21</u> <u>thresh = 0.001</u> <u>fast.k = 20</u> <u>pmethod = backward</u>	degree = {1, 2, 3} nprune = {5, 10, 15, 20, 25}
ANN	nnet	<u>weights = 1</u> <u>rang = 0.7</u> <u>maxit = 100</u> <u>maxNWts = 1000</u> <u>abstol = 10⁻⁴</u> <u>reltol = 10⁻⁸</u>	size = {1, 2, 4, 8, 20} decay = {0.0, 0.1, 0.5, 1.0, 2.0}
RF	randomForest	<u>ntree = 500</u> <u>sampsize = 528</u> <u>nodesize = 5</u> <u>nPerm = 1</u>	mtry = {2, 3, 4, 5, 6, 7}
M5	Cubist	<u>rules = 100</u> <u>extrapolation = 100</u> <u>sample = 0</u>	committees = {10, 50, 100} neighbors = {0, 5, 9}

2

1

2 Table 3. Cross validation errors for each assessed model.

Standard Formulation		GLM	GAM	MARS	RF	M5	ANN	Climatology	Noah LSM
MAE	Gilgel Abbay	30.78	18.54	16.75	14.89	15.11	17.22	10.42	28.11
	Gumara	4.29	3.41	3.28	2.67	2.96	3.15	2.57	3.95
	Koga	1.50	1.30	1.38	1.20	1.17	1.23	1.06	1.97
	Megech	4.45	2.64	2.83	2.37	2.53	3.04	2.54	4.09
	Ribb	4.69	2.98	3.50	2.97	3.27	3.17	2.81	7.01
NSE	Gilgel Abbay	-0.02	0.81	0.83	0.87	0.86	0.84	0.95	0.59
	Gumara	0.04	0.51	0.61	0.80	0.66	0.70	0.81	0.48
	Koga	0.45	0.71	0.65	0.76	0.77	0.76	0.83	0.25
	Megech	-1.85	0.63	0.46	0.73	0.65	0.52	0.71	0.41
	Ribb	-1.14	0.71	0.39	0.71	0.31	0.67	0.73	-0.75
Anomaly Formulation		GLM	GAM	MARS	RF	M5	ANN	Climatology	Noah LSM
MAE	Gilgel Abbay	9.73	9.82	10.10	10.12	9.94	9.79	10.42	28.11
	Gumara	2.22	2.25	2.43	2.23	2.16	2.22	2.57	3.95
	Koga	1.03	1.06	1.08	1.09	1.05	1.05	1.06	1.97
	Megech	2.49	2.48	2.63	2.66	2.69	2.50	2.54	4.09
	Ribb	2.79	2.76	2.84	2.70	2.78	2.77	2.81	7.01
NSE	Gilgel Abbay	0.95	0.95	0.95	0.95	0.95	0.95	0.95	0.59
	Gumara	0.85	0.85	0.82	0.85	0.86	0.86	0.81	0.48
	Koga	0.83	0.82	0.81	0.81	0.82	0.82	0.83	0.25
	Megech	0.73	0.72	0.65	0.66	0.61	0.72	0.71	0.41
	Ribb	0.73	0.75	0.72	0.75	0.73	0.74	0.73	-0.75

3

1

2 Table 4. Residual autocorrelation factors at a 12-month lag for the standard formulation and
3 anomaly formulation models, and resulting mean annual observed and predicted flow.

4

	Autocorrelation Factors		Mean Annual Flow (MCM)		
	Standard	Anomaly	Observed	Standard	Anomaly
Gilgel	0.33	0.11	22,925	20,703	22,958
Gumara	0.29	0.07	2,870	2,392	2,734
Koga	0.04	0.10	1,383	1,333	1,386
Megech	0.05	0.04	2,035	1,637	2,028
Ribb	0.21	-0.01	2,575	1,969	2,615

5

1 Table 5. Covariate importance measurements from each basin's model

Model type	Linear model				Generalized additive model		M5 model tree		Random forest
Measure of influence	Linear regression coefficients and associated p-values				Estimated degrees of freedom (EDF) and associated p-values		Covariate usage in tree rules and model coefficients		Increase in MSE when covariate is randomly permuted
Basin	Gilgel Abbay		Koga		Megech		Gumara		Ribb
Covariate	Coefficient estimate	P-value	Coefficient estimate	P-value	EDF	P-value	Tree rules	Model coefficients	Percent increase in MSE
Prec	0.22	< 0.01	0.24	< 0.01	1.346	< 0.01	5%	58%	7.71%
Prec (lag 1)	0.10	0.03	0.16	< 0.01	0.624	0.08	0%	19%	2.79%
Prec (lag 2)	0.01	0.74	0.05	0.26	0	0.29	0%	0%	1.10%
Temp	-0.09	0.08	-0.07	0.17	1.023	0.07	0%	47%	12.74%
Temp (lag 1)	-0.04	0.49	-0.06	0.22	0	0.32	0%	46%	4.97%
Temp (lag 2)	-0.01	0.81	-0.09	0.08	0	0.56	0%	0%	8.16%
Agr. LC	0.00	0.33	0.02	0.01	1.986	< 0.01	86%	73%	15.21%

2

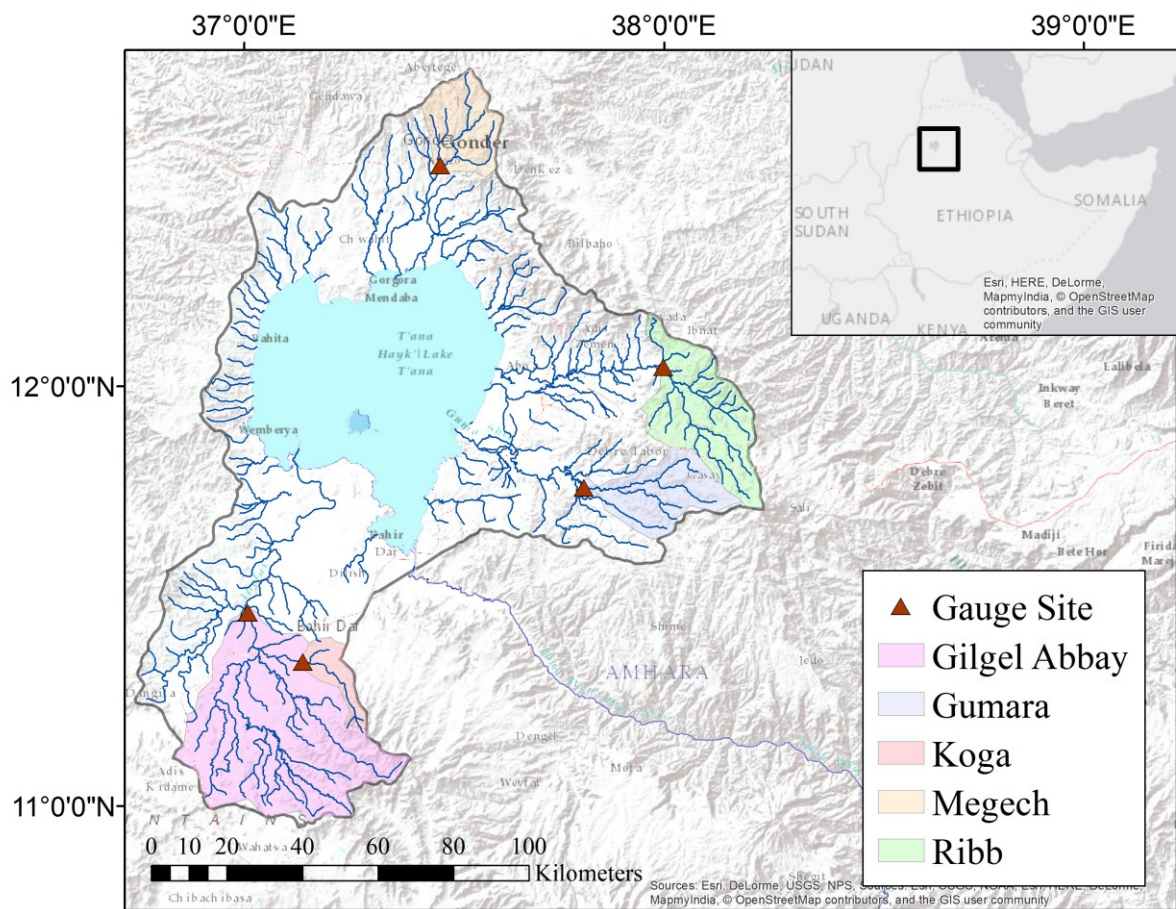
1 Table 6. Mean and standard deviation values for temperature, wet-season rainfall, and dry-
 2 season rainfall in each basin.

3

	Temperature (°C)		Wet season rainfall (mm/month)		Dry season rainfall (mm/month)	
	Mean	SD	Mean	SD	Mean	SD
Gilgel Abbay	15.7	1.54	206	145	39.3	56.5
Gumara	17.7	1.55	186	137	29.0	43.6
Koga	15.7	1.54	206	145	39.3	56.5
Megech	20.6	1.75	234	118	41.4	60.9
Ribb	18.2	1.61	263	115	45.8	57.0

4

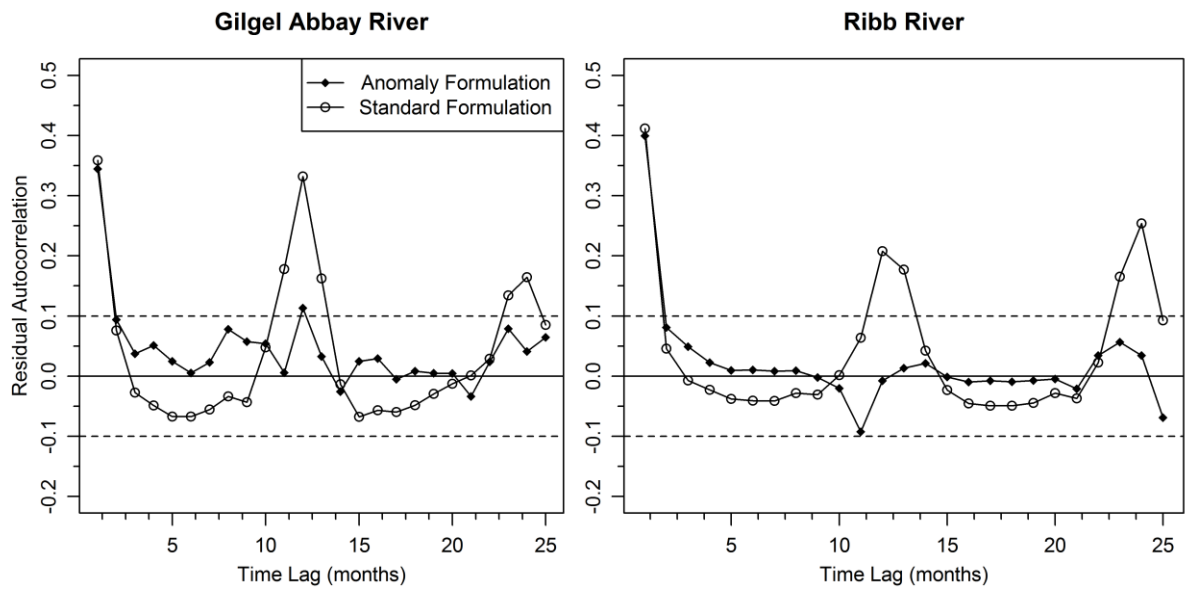
1 Figure 1. Map of Lake Tana and surrounding rivers



2

3

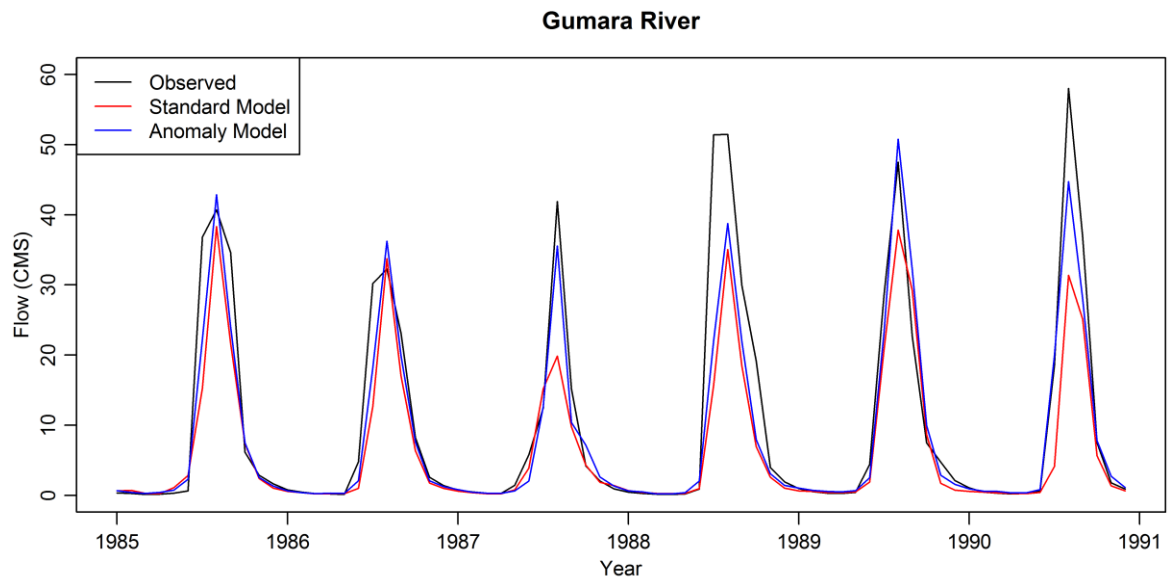
1 Figure 2. Autocorrelation in model residuals for the Gilgel Abbay and Ribb Rivers



2

3

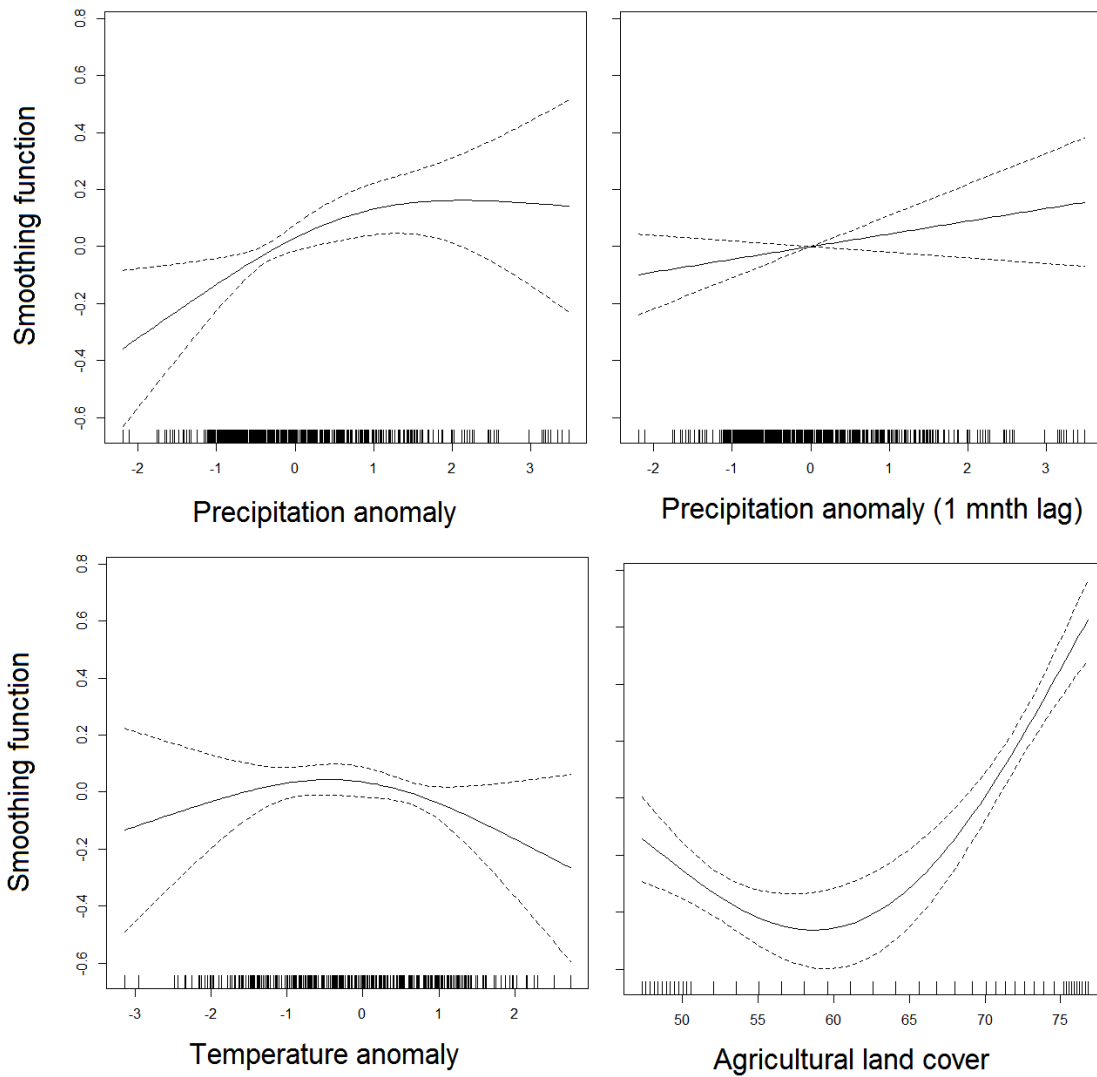
1 Figure 3. Example observed and predicted flows for Gilgel Abbay River from 1995 to 2000.



2

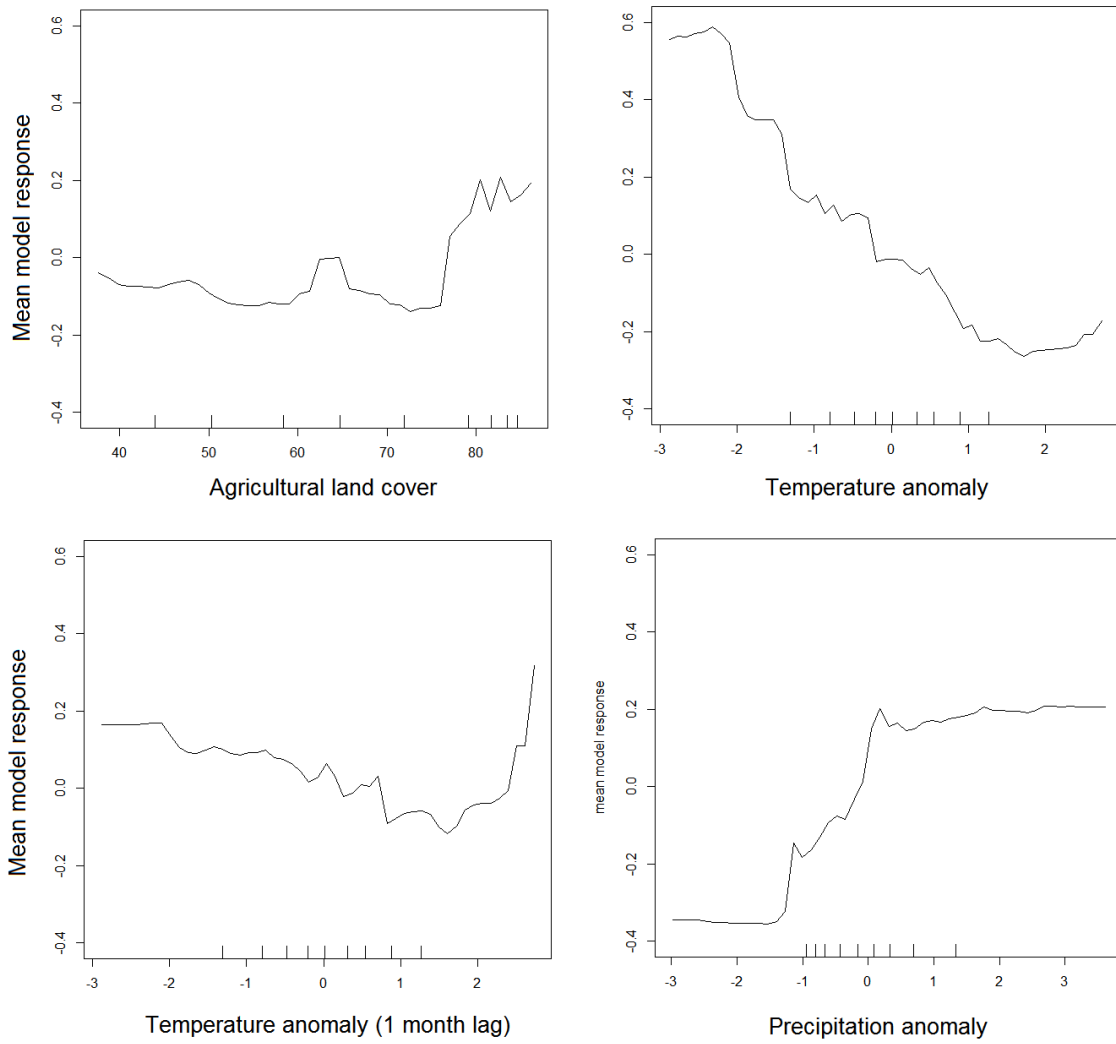
3

1 Figure 4. Plots of the smoothing functions used in the Megech River GAM. Hash marks along
2 the x-axis indicate observation values of each covariate.



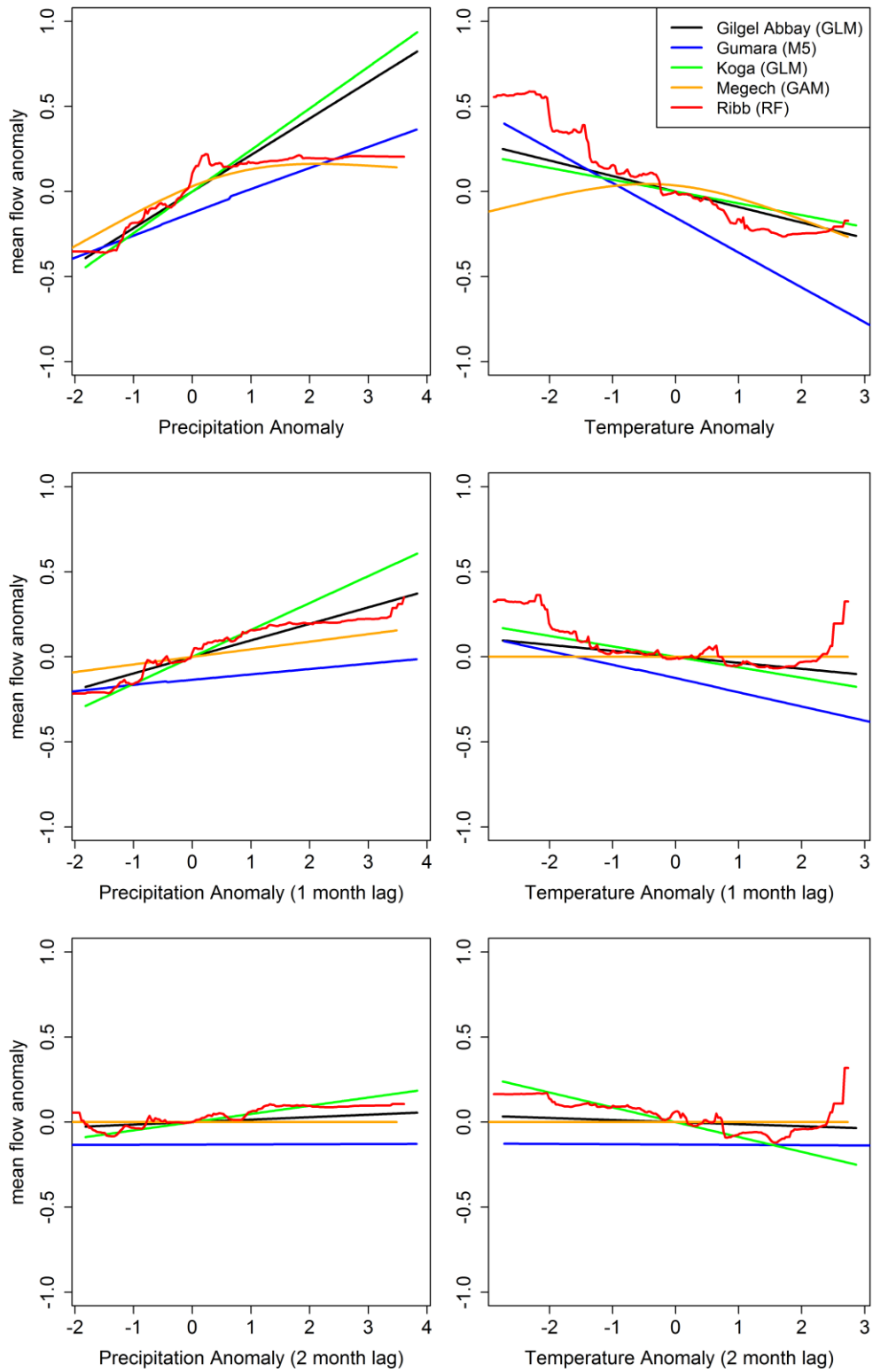
3
4

1 Figure 5. Partial dependence plots for the Ribb River random forest model. Hash marks along
2 the x-axis show covariate sample decile values.



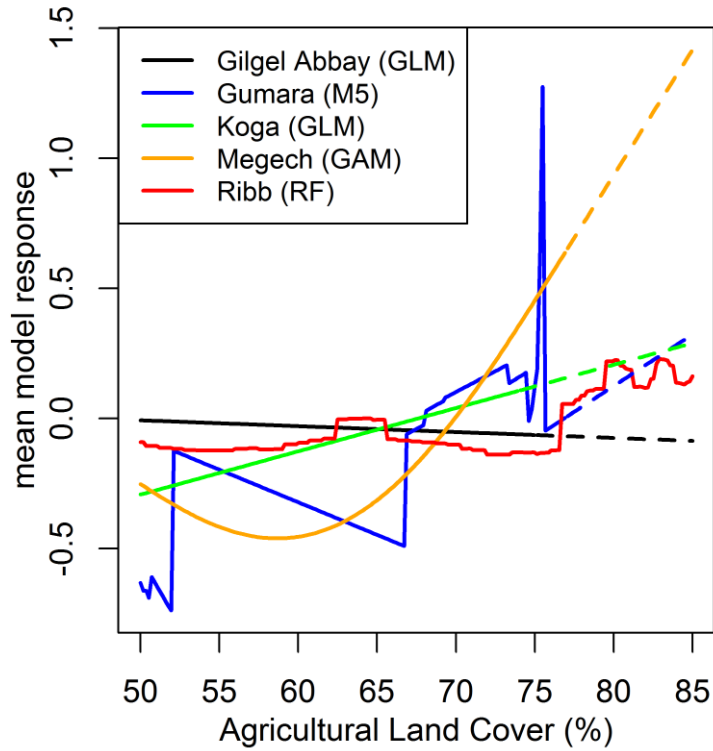
3
4

1 Figure 6. Partial dependence plots for climate covariates in the highest performing model in
2 each basin. Model type is indicated in parentheses.



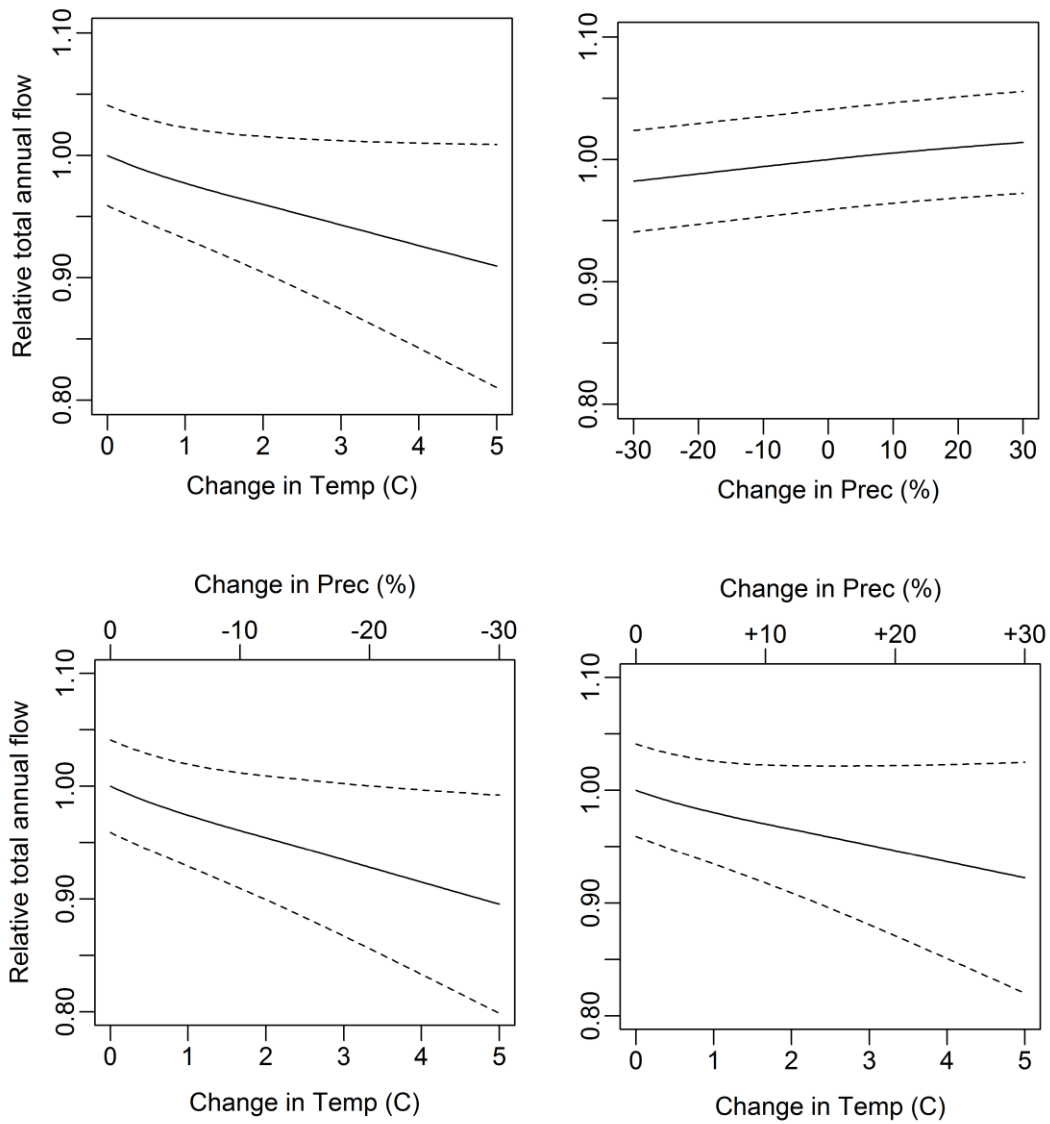
3

1 Figure 7. Partial dependence plot for agricultural land cover in the highest performing model
2 in each basin. Model type is listed in parentheses for each basin. Dashed lines
3 indicate values that exceed historic levels of agricultural land cover experienced in
4 that basin.



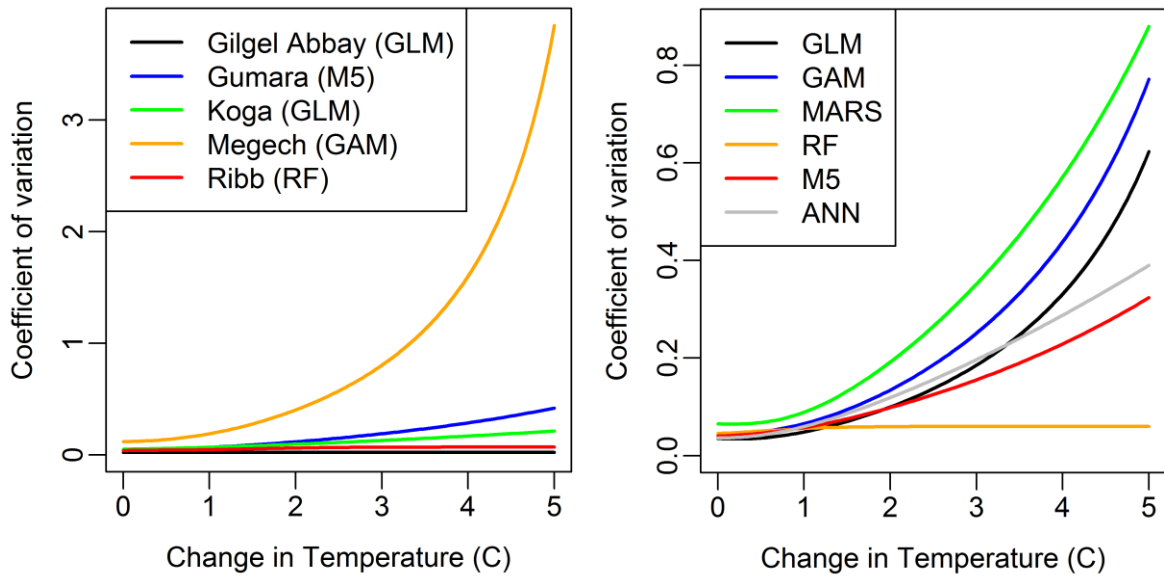
5
6

1 Figure 8. Projected changes in total streamflow (relative to current long-term average) under
 2 changing climate conditions. The top two panels show the sensitivity to changes in
 3 temperature and precipitation when they are varied independently. The bottom panel
 4 shows sensitivity to changing temperature in conjunction with decreasing (left
 5 panel) and increasing (right panel) precipitation. Dashed lines represent 95%
 6 confidence bounds from bootstrap resampling.



7
8

1 Figure 9. Changes in the coefficient of variation across bootstrap resamples from the highest
2 performing model in each basin (left panel) and multiple models all applied to the
3 Gumara basin (right panel).



4