

1 **Machine learning methods for empirical streamflow**  
2 **simulation: a comparison of model accuracy,**  
3 **interpretability and uncertainty in seasonal watersheds**

4

5 **J. E. Shortridge,<sup>1</sup> S. D. Guikema,<sup>2</sup> and B. F. Zaitchik<sup>3</sup>**

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

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

10 [3]{Department of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, USA}

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

12

13

1 **Abstract**

2 In the past decade, machine-learning methods for empirical rainfall-runoff modeling have  
3 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 often evaluated model  
8 performance based on predictive accuracy alone, while not considering broader objectives  
9 such as model interpretability and uncertainty that are important if such methods are to be  
10 used for planning and management decisions. In this paper, we use multiple regression and  
11 machine-learning approaches (including generalized additive models, multivariate adaptive  
12 regression splines, artificial neural networks, random forests, and M5 cubist models) to  
13 simulate monthly streamflow in five highly-seasonal rivers in the highlands of Ethiopia and  
14 compare their performance in terms of predictive accuracy, error structure and bias, model  
15 interpretability, and uncertainty when faced with extreme climate conditions. While the  
16 relative predictive performance of models differed across basins, data-driven approaches were  
17 able to achieve reduced errors when compared to physical models developed for the region.  
18 Methods such as random forests and generalized additive models may have advantages in  
19 terms of visualization and interpretation of model structure, which can be useful in providing  
20 insights into physical watershed function. However, the uncertainty associated with model  
21 predictions under extreme climate conditions should be carefully evaluated, since certain  
22 models (especially generalized additive models and multivariate adaptive regression splines)  
23 become 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 machines  
8   have been shown to be powerful tools for predictive modeling and exploratory data analysis,  
9   particularly in systems that exhibit complex, non-linear behavior (Solomatine and Ostfield,  
10   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 While a number of comparison studies exist that apply multiple empirical models to a  
27 given problem, finding generalizable insights from these studies is hindered because of the  
28 limited number of models and datasets evaluated. Perhaps the most comprehensive  
29 comparison to date is that of Elshorbagy et al. (2010a and 2010b), who compared six methods  
30 for data-driven modeling of daily discharge in the Ourthe River in Belgium. This work found  
31 that linear models were able to perform comparably to much more complex methods when the  
32 data content of the models were limited, or when system input-output behavior was close to

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

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

21       In this work, we compare six methods for empirical streamflow simulation (linear  
22 models, generalized additive models, multivariate adaptive regression splines, random forests,  
23 M5 model trees and ANNs) in five rivers in the Lake Tana basin in Ethiopia. This study  
24 region was selected as it provides insights into the use of data-driven models for streamflow  
25 simulation in tropical regions of the world that are underrepresented in existing studies; for  
26 instance, a review of 210 articles on water resource applications of ANNs found that over  
27 three quarters of the studies evaluated were conducted in North America, Europe, Australia,  
28 or temperate East Asia (Maier et al., 2010). Existing studies conducted in tropical regions  
29 generally apply a single methodology to the basin of interest and evaluate predictive accuracy  
30 alone (see for instance, Machado et al., 2011; Chibanga et al., 2003; Antar et al., 2006; Aqil et  
31 al., 2007), making it difficult to find generalizable insights into the relative advantages of  
32 different modeling approaches in these regions. Better development of data-driven models for

1 these regions has the potential to be particularly valuable because data limitations and  
2 complex hydrodynamic processes often hinder the use of physical watershed models, but  
3 relatively long time series of streamflow, precipitation and temperature may be available at a  
4 monthly timescale. These data, combined with information on relevant landscape change (in  
5 particular, the expansion of agricultural land cover), can be leveraged to create reasonably  
6 accurate empirical models.

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

## 17 **2 Data and Methods**

### 18 **2.1 Study Area**

19 Lake Tana is located at an elevation of approximately 1800 meters in the highlands of  
20 northwest Ethiopia (Fig. 1). The catchment draining to the lake encompasses approximately  
21 12,000 square kilometers, and the four main tributaries providing water to the lake are the  
22 Gilgel Abbay (including its tributary, the Koga River), Ribb, Gumara, and Megech Rivers.  
23 Collectively, these rivers account for 93% of the inflow to the lake (Alemayehu et al., 2010).  
24 Ninety percent of rainfall in the basin occurs during the wet season from May until October,  
25 and there is significant interannual variability in precipitation with annual rainfall levels  
26 ranging from below 1000 mm to over 1800 mm (Achenef et al., 2013). Population growth and  
27 expansion of agricultural and pastoral land use in the region has resulted in substantial  
28 deforestation and land degradation, with agricultural, pastoral and settled land cover  
29 comprising over 70% of the basin’s surface area (Rientjes et al., 2011; Garede and Minale,  
30 2014; Gebrehiwot et al., 2010). There is some evidence that this has impacted the hydrology

1 of the rivers draining into the lake (Gebrehiwot et al., 2010). A summary of basin  
2 characteristics for the evaluation period of 1960-2004 is presented in Table 1.

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

## 30 **2.2 Data and Model Development**

31        Models were developed using monthly streamflow, climate, and land cover data for  
32 the period from 1961 to 2004, resulting in 528 monthly observations. In each of the five major

1 rivers in the basin, we developed empirical models that estimated monthly streamflow as a  
 2 function of climate conditions and agricultural land cover in each basin. Monthly streamflow  
 3 data were taken from historic stream gauge records for each basin, as reported in feasibility  
 4 studies developed for proposed irrigation projects (Alemayehu, 2010). Historic data for  
 5 monthly average temperature and monthly total precipitation in each river basin were derived  
 6 from the University of East Anglia Climate Research Unit (CRU) TS3.10 gridded  
 7 meteorological fields (Harris et al., 2014), which are based on meteorological station  
 8 observations. Finally, to account for historic increases in agricultural and pastoral land cover  
 9 that have occurred in the basin, the percentage of land cover used for any crop or grazing was  
 10 estimated from historic land cover analyses described by Rientjes et al. (2011), Gebrehiwot et  
 11 al. (2010), and Garede and Minale (2014). These studies used historic aerial photos and  
 12 satellite images to estimate land cover changes in the Ribb, Gilgel Abbay, and Koga basins  
 13 from the periods of 1957 to 2011. The percentage of agricultural land cover was interpolated  
 14 for years when data weren't available, and the value of agricultural land cover in the two  
 15 basins without data was assumed to be equal to average agricultural land cover in the basins  
 16 with data. Land cover was assumed to change on an annual, rather than monthly basis. While  
 17 this approach is prone to errors that could stem from differing rates of land use change  
 18 through time and between basins, it does provide a mechanism for capturing the long-term  
 19 trend of expanding agricultural land cover that has been observed throughout the Ethiopian  
 20 highlands when detailed land-cover data are unavailable. Including this data improved out-of-  
 21 sample predictive accuracy of the models, further suggesting that it was a valuable addition.

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

$$24 \quad \log(Q_{b,t}) = f(P_{b,t}, P_{b,t-1}, P_{b,t-2}, T_{b,t}, T_{b,t-1}, T_{b,t-2}, AgLC_{b,t}) + \varepsilon_{b,t} \quad (1)$$

25 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  
 26 total precipitation and average temperature in river basin  $b$  at time period  $t$ ,  $AgLC_{b,t}$  is the total  
 27 percentage of agricultural land cover in basin  $b$  at time  $t$ , and  $\varepsilon_{b,t}$  is the model error. The  
 28 subscripts  $t-1$  and  $t-2$  indicate lagged measurements from one and two months prior, and were  
 29 included to roughly account for storage times longer than one month that could impact  
 30 streamflow in each river. While the exact time of concentration is not known in each basin,  
 31 the minor influence of climate conditions at two months prior suggest that climate conditions



1 from beyond this time period do not contribute significantly to flow variability. The function  $f$   
 2 represents a general function that differed between the specific models assessed and is  
 3 discussed in more detail below. The logarithm of monthly streamflow was used as a response  
 4 variable to keep model predictions positive. The distribution of streamflow data and log-  
 5 transformed streamflow values in each basin are shown in supplementary Fig. S1.

6 In the second formulation, streamflow and climate anomalies were used as the  
 7 response and predictor variables to better account for the highly seasonal nature of streamflow  
 8 and precipitation in the region. Streamflow anomalies were calculated for each observation by  
 9 subtracting the long-term average streamflow for that month ( $m$ ) from the observed value and  
 10 dividing this number by the long-term standard deviation of that month's streamflow as in Eq.  
 11 (2). Anomaly values thus represent how streamflow in a given month compares to the long-  
 12 term average flow for that month; for instance, an anomaly value of 1.0 for June of 1990  
 13 would indicate that streamflow in that month was one standard deviation higher than the  
 14 average June flow from 1961 to 2004. This procedure was repeated for precipitation and  
 15 temperature, and these values were then used to fit models of the form described in Eq. (3). In  
 16 each month of the time series, the model estimates the flow relative to the long-term average  
 17 flow for that month, based on whether temperature and precipitation values were greater or  
 18 less than their long-term averages, as well as the percentage of agricultural land cover in that  
 19 month of the time series. In this sense, the anomaly values are calculated based on climatic  
 20 and land cover conditions that vary through time. These anomaly values are then converted  
 21 back to raw flow values based on the long-term average and standard deviation of flow for  
 22 that month. The distribution of streamflow anomaly values in each basin are shown in  
 23 supplementary Fig. S1.

$$24 \quad Q_{b,t}^{AN} = \frac{Q_{b,t} - \bar{Q}_{b,m}}{sd(Q_{b,m})} \quad (2)$$

$$26 \quad Q_{b,t}^{AN} = f(P_{b,t}^{AN}, P_{b,t-1}^{AN}, P_{b,t-2}^{AN}, T_{b,t}^{AN}, T_{b,t-1}^{AN}, T_{b,t-2}^{AN}, AgLC_{b,t}) + \varepsilon_{b,t} \quad (3)$$

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

- 28 1. A Gaussian linear regression model (GLM) using the basic stats package in the R  
 29 statistical computing software (R Development Core Team, 2014)
- 30 2. Gaussian generalized additive model (GAM): GAMs are a semi-parametric  
 31 regression approach where the response variable is estimated as the sum of

1 smoothing functions applied over predictor variables. These functions allow the  
2 model to capture non-linear relationships between the predictor and response  
3 variables without *a priori* assumptions about the form (eg., quadratic, logarithmic)  
4 of these functions, and are fit using penalized likelihood maximization to prevent  
5 model overfitting (Hastie and Tibshirani, 1990). GAMs were fit using the *mgcv*  
6 package in R (Wood, 2011).

- 7 3. Multivariate adaptive regression splines (MARS): MARS are a non-parametric  
8 regression approach where the response variable is estimated as the sum of basis  
9 functions fit to recursively partitioned segments of the data (Friedman, 1991).  
10 MARS models were fit using the *earth* package in R (Milborrow, 2015).
- 11 4. Artificial neural network (ANN): ANNs are a non-parametric regression approach  
12 represented by a network of nodes and links that connects predictor variables to  
13 the response variable. Each link in the network represents a function that maps the  
14 input nodes into the output node (Ripley, 1996). ANN models were fit using the  
15 *nnet* package in R (Venables and Ripley, 2013).
- 16 5. Random forest (RF): Random forests are a rule-based, non-parametric regression  
17 approach where the model prediction is created by averaging the predicted value  
18 from multiple regression trees which are trained on separate bootstrapped  
19 resamples of the data. Each tree is fit using a small, randomly selected subset of  
20 predictor variables, resulting in reduced correlation between trees (Breiman,  
21 2001). Random forest models were fit using the *randomForest* package in R (Liaw  
22 and Wiener, 2002).
- 23 6. M5 model: M5 models are a rule-based, non-parametric regression approach that  
24 fits a linear regression model to each terminal node of a regression tree (Quinlan,  
25 1992). M5 models were fit using the *Cubist* package in R (Kuhn et al., 2014).
- 26 7. Climatology model: A climatology model that simply predicted each month's  
27 streamflow as equivalent to the long-term average streamflow for that month was  
28 included for comparison purposes.

### 29 **2.3 Model Evaluation**

30 When using non-parametric regression approaches, it is important to avoid overfitting a  
31 model to a given dataset because this can result in large errors in out-of-sample predictions  
32 (Hastie et al., 2009). To avoid model overfit, the *caret* package in R (Kuhn, 2015) was used to

1 determine model parameters for the MARS, ANN, RF and M5 models. This package uses  
2 resampling to evaluate the effect that model parameters have on the model's predictive  
3 performance and chooses the set of parameters that minimizes out-of-sample error (Kuhn  
4 2015). In this evaluation, 25 bootstrap resamples of the training dataset were generated for  
5 each parameter value to be assessed. A model was fit using each bootstrap sample and used to  
6 predict the remaining observations, and the parameter values that minimized average RMSE  
7 across all resamples. Details on the specific parameters evaluated for each model are  
8 presented in Table 2. While the development of more complex structures are possible for  
9 some models, this process can result in over-parameterization and poor model performance  
10 (Gaume and Gosset, 2003; Han et al., 2007). Additionally, the use of a standardized  
11 parameterization procedure allows for a more even comparison between different model  
12 types.

13 The predictive ability of each model was assessed using 50 random holdout cross-  
14 validation samples. In each sample, a random selection of years were chosen, and  
15 observations from these years were removed ("held-out") from the dataset. The size of the  
16 held-out sample ranged from 1 to 9 years. Each model was then fit to the remaining portion of  
17 the data, using the caret package described above to determine model parameters for the  
18 MARS, ANN, RF and M5 models. These models were then used to predict streamflow for the  
19 held-out portion of the data, and both the mean absolute error (MAE) and NSE were  
20 calculated after transforming model predictions after back to the original streamflow units.  
21 Mean MAE and NSE were calculated for each model across the 50 cross-validation samples  
22 and used to choose the model with the highest predictive accuracy in each basin. This cross-  
23 validation procedure provides a mechanism for evaluating how well a model will generalize  
24 to an unseen set of data while avoiding some of the problems that can arise from the use of a  
25 single calibration and validation dataset (Elshorbagy et al., 2010a; Han et al., 2007).

26 MAE was included as an error metric because it provides a simple and easily  
27 interpretable measure of error on the same scale as observed flow volumes. While NSE values  
28 are acknowledged to be a flawed performance metric in highly seasonal watersheds where  
29 seasonal fluctuations contribute to a substantial portion of flow variability (Schaeffli and  
30 Gupta, 2007; Legates and McCabe, 1999), this metric was included to provide a rough  
31 comparison of how empirical model performance compared to the performance of physical  
32 models developed for the region. The use of alternative error metrics has been discussed

1 extensively in the literature (for instance Pushpalatha et al., 2012; Mathevet et al., 2006; Criss  
2 and Winston, 2008), and could provide additional insights into what contributes to predictive  
3 capabilities of different model formulations. However, this work examined predictive  
4 accuracy based on MAE and NSE alone to allow for greater focus on how models differ in  
5 terms of error structure and uncertainty.

6 As a rough point of comparison for the statistical models developed in this research, we  
7 also evaluated discharge estimates derived from a process-based hydrological model. The  
8 model used in this application is the Noah Land Surface Model version 3.2 (Noah LSM; Ek  
9 et. al, 2003; Chen et al., 1996). Noah LSM was implemented for offline simulations of the  
10 Lake Tana basin at a gridded spatial resolution of 5km for the period 1979-2010 using a time  
11 step of 30 minutes. Meteorological forcing was drawn from the Princeton 50-year reanalysis  
12 dataset (Sheffield et al. 2006), downscaled to account for Ethiopia's steep terrain using  
13 MicroMet elevation correction equations (Liston & Elder 2006). The Princeton reanalysis was  
14 selected because it provides relatively high resolution meteorological fields, including all  
15 variables required to run a water and energy balance LSM like Noah, for the period 1948-  
16 present. While higher resolution and possibly higher quality datasets are available for recent  
17 years, this longer dataset was utilized to compare the process-based model to statistical  
18 models developed for a long historical period. Soil parameters for the Noah simulation were  
19 drawn from the FAO global soil database, land use was defined according to the United States  
20 Geological Survey (USGS) global 1km land cover product, and vegetation fraction was  
21 derived from MODerate Imaging Spectroradiometer (MODIS) imagery. Land cover was  
22 treated as a static parameter over the full length of the simulation, as spatially complete  
23 estimates of historical land use were not available at the required resolution and specificity.

24 The highest performing model in each basin based on MAE was retained for more  
25 detailed evaluation of model error structure, covariate influence, and uncertainty in climate  
26 change sensitivity analysis. To generate a complete time-series of out-of-sample model  
27 predictions for error analysis, the holdout cross validation procedure was repeated for the  
28 highest performing standard-formulation and anomaly-formulation models for each basin, but  
29 this time holding out a single year of observations in each iteration. The predictions from this  
30 cross validation were used to evaluate the how model error structure might impact model  
31 predictions used for water resource applications. The influence of different predictor variables  
32 on model predictions was also assessed for the highest performing model in each basin after

1 being fit to the complete dataset. Each predictor variable was assessed using metrics for  
2 covariate importance and influence that are unique to that model type, demonstrating how  
3 models could be used to gain physical insights about data-scarce regions and the mechanisms  
4 for generating these insights for each type of model. Partial dependence plots (Hastie et al.,  
5 2009) were also generated for each covariate for the highest performing model in each basin  
6 to provide insights about how covariate influence compared across different basins and model  
7 types.

8 Finally, two evaluations were conducted to assess uncertainty in model projections of  
9 streamflow under increasingly extreme climate conditions to better understand the  
10 implications of using different model formulations for climate change impact studies. Model  
11 projections of streamflow in different climate conditions are likely to be accompanied by  
12 considerable uncertainty, particularly when climate conditions exceed those experienced  
13 historically. To assess this uncertainty, the best performing model in each basin was used to  
14 generate streamflow predictions for 1) changes in temperature from 0 to 5° C, 2) changes in  
15 precipitation from -30 to +30%, 3) an increase in temperature to 5° C combined with a  
16 decrease in precipitation to -30%, and 4) an increase in temperature to 5° C combined with an  
17 increase in precipitation to +30%. For each of the four assessments, the models generated  
18 predictions for the 45-year historic climate record adjusted for a given degree of climate  
19 change using the delta-change method (Gleick, 1986), while holding agricultural land cover  
20 constant at 60%. In this method, monthly temperature values are simply added to the  
21 temperature change value, and monthly precipitation values are multiplied by the precipitation  
22 change percentage. Model predictions for the altered climate record were then used to  
23 calculate the average annual streamflow in each river. This process was repeated 100 times  
24 for models fit on random bootstrap resamples of the historic dataset to generate uncertainty  
25 bounds surrounding model predictions and evaluated how the uncertainty in these predictions  
26 increased as climate conditions became more extreme. It is important to recognize that these  
27 should not be interpreted as a prediction or assessment of actual climate change impacts, but  
28 rather a measurement of the sensitivity of modeled streamflow in the basin to different  
29 climate conditions. Since one of the key motivations for using rainfall-runoff models is to  
30 understand how climate change may impact water resources, it is important to understand  
31 how model formulation contributes to this sensitivity and uncertainty.

## 1 **3 Results**

### 2 **3.1 Model Accuracy and Error Structure**

3 Table 3 shows the out-of-sample cross validation errors for each model assessed in each  
4 basin. The random forest model had the lowest mean absolute error for the standard-  
5 formulation model in four of the five basins, with the M5 model performing best for the Koga  
6 basin. These models outperformed the Noah LSM simulations in all basins assessed. The  
7 Noah LSM errors are for a single period of analysis and thus don't present an exact corollary  
8 to the cross validation performed for the empirical models. Nevertheless, the significant  
9 increases in errors associated with the Noah LSM model demonstrates the difficulty  
10 associated with the use of process-based models in the region, particularly when relying on  
11 global datasets that may be unreliable at the spatial and temporal resolutions required for  
12 physical modeling. Physical models developed for monthly streamflow prediction in other  
13 basins within the Ethiopian highlands have reported NSE values ranging from 0.53 to 0.92  
14 (van Griensven et al., 2012), compared to values ranging from 0.71 to 0.87 for the random  
15 forest models developed here. If this measure alone was used for model evaluation, these  
16 empirical models would generally be classified as having good performance based on the  
17 guidelines suggested by Moraisi et al. (2007). However, the climatology model outperforms  
18 the best standard formulation models in all basins except Megech, indicating that in the  
19 majority of basins the errors from the fitted empirical models are higher than those that result  
20 from simply using the long-term monthly average for each month's prediction. This is due to  
21 the fact that seasonality accounts for such a large portion of the variability in monthly flow  
22 values, and demonstrates how high NSE values can be quite easy to obtain in seasonal basins.

23 Evaluation of anomaly model errors indicates that the models using this formulation  
24 achieve better predictive accuracy than those using the standard formulation, and are able to  
25 outperform the climatology model based on both NSE and MAE in all basins. However, the  
26 highest performing models in each basin varies more when the anomaly formulation is used,  
27 with the GLM, GAM, random forest, and M5 models all minimizing MAE in different basins.  
28 In all basins except Koga, the highest performing model significantly outperformed the  
29 climatology model based on paired Wilcoxon rank-sum tests (Bonferroni-corrected p-value <  
30 0.01).

1 Further exploration of model residuals indicates another important advantage of using  
2 the anomaly model formulation. In the standard model formulation, model residuals appear to  
3 be non-random. Example autocorrelation plots are shown for the Gilgel Abbay and Ribb  
4 Rivers in Fig. 2, and demonstrate that a positive autocorrelation exists at the 12 month time  
5 lag. For brevity, only plots for two rivers are shown, although this autocorrelation existed in  
6 the standard-formulation models for all basins except Megech (Table 4). This autocorrelation  
7 occurs because the standard-formulation models consistently underestimate wet-season  
8 streamflow while overestimating dry-season flows, as is apparent in hydrographs of observed  
9 and predicted streamflow (Fig. 3). Because wet-season flows contribute such a large portion  
10 of the total annual flow volume, this results in regular underestimation of aggregate values  
11 such as mean annual flow (Table 4). This autocorrelation is reduced in the anomaly-  
12 formulation models, meaning that they are better able to capture the peak flow volumes  
13 experienced in the wet season and do not underestimate mean annual flow to the same degree  
14 that the standard formulation models do.

### 15 **3.2 Model Structure and Covariate Influence**

16 Evaluating the relationship between predictor covariates and streamflow response can  
17 lend insight into the physical processes underlying runoff generation in each basin. There are  
18 two components of this relationship that can be evaluated: how much each covariate  
19 contributes to model accuracy (covariate importance), and the direction and nature of the  
20 relationship between covariate values and model response (covariate influence). In many  
21 machine-learning models, complete description of the all of the mathematical relationships  
22 within the model (for instance, through description of each tree comprising a random forest  
23 model) is infeasible, requiring the use of other mechanisms for understanding covariate  
24 importance and influence. However, because each model type is structured in a different way,  
25 these mechanisms differ. This section first describes the mechanisms available for obtaining  
26 insights about covariate influence in each of the highest performing models. To provide a  
27 mechanism for comparing results across different basins, each basin model is then assessed  
28 using the general approach of partial dependence plots.

29 In the Gilgel Abbay and Koga basins, the highest performing model was a simple  
30 linear regression model. These models can be evaluated by reviewing model coefficients and  
31 associated p-values, as shown in Table 5. In a standard linear regression, model coefficients  
32 can be interpreted as the mean change in the response variable that results from a unit change

1 in that covariate when all others are held constant. These coefficients are for streamflow  
2 anomalies rather than raw values, making their immediate interpretation less intuitive. For  
3 instance, in the Gilgel Abbay model an increase of one standard deviation in precipitation  
4 results in an increase of 0.22 standard deviations in flow. The associated p-value for each  
5 coefficient evaluates a null hypothesis that the true coefficient value is equal to zero given the  
6 other covariates in the model, and thus has no influence on the response variable.

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

15 Interpretation of variable influence in GAMs is based on the estimated degrees of  
16 freedom (EDF) a covariate's smoothing function  $s(X_i)$  uses within a model (Hastie and  
17 Tibushini, 1986). An EDF value of one or below indicates a linear function relating the  
18 response variable to that covariate, while values greater than one represent a non-linear  
19 smoothing function. An EDF value of zero indicates that the covariate smoothing function is  
20 penalized to zero (meaning it has no influence on model predictions). In the model for the  
21 Megech River, the terms for lagged temperature at one and two months, as well as  
22 precipitation lagged at two months were all smoothed to zero. Of the remaining covariates,  
23 lagged precipitation has a linear impact on model response, while precipitation, temperature  
24 and land cover have non-linear impacts. Smoothing functions can be plotted to gain more  
25 insight about these relationships (Fig. 4). The functions for precipitation anomaly, lagged (one  
26 month) precipitation anomaly, and agricultural land cover show a positive relationships with  
27 streamflow, while the function for temperature anomaly predicts low streamflow at both high  
28 and low anomalies.

29 P-values test the null hypothesis that a covariate's smoothing function is equal to zero,  
30 but rest on the assumption that model residuals are homoscedastic and independent (Wood,  
31 2012). Similar to the linear models, residuals in the Megech GAM model appear to be both  
32 autocorrelated and heteroscedastic, meaning that a formal statistical interpretation of this



1 value may be inappropriate and that confidence bounds around smoothing functions might be  
2 misleading.

3 The M5 cubist model fit for the Gumara basin is an ensemble of 100 small M5  
4 regression trees. In each tree, the model splits observations based on logical rules related to  
5 one or more covariates and fits a linear regression model to each set of observations. The final  
6 model prediction is the average across all of the individual trees. Using this sort of ensemble  
7 approach can reduce model variance and improve accuracy if the individual trees are  
8 unbiased, uncorrelated predictors (Breiman 1996). This can be useful in avoiding models that  
9 are overfit to the data, but can reduce model interpretability since direct visualization of  
10 model structure becomes impractical as the number of trees increases. However, the  
11 frequency with which individual covariates are used as splitting points within trees and as  
12 regression coefficients can provide some insights about covariate importance (Table 5; note  
13 that because multiple covariates can be used for rules and linear models, these don't  
14 necessarily add to 100%). Model rules were largely based on land cover, with some rules  
15 based on precipitation. These two covariates were also used most frequently in linear  
16 regressions at model nodes, followed by temperature (current and 1-month lag) and 1-month  
17 lagged precipitation. Notably, climate data from 2 months lagged were not used at all. While  
18 this can be useful in identifying which covariates have the largest impact on model  
19 predictions, it doesn't provide any information regarding the nature or direction of that  
20 influence.

21 Similarly, the random forest model developed for the Ribb basin is an ensemble of  
22 regression trees in which the final model prediction is the average of the predictions from  
23 each individual tree. However, random forests use standard regression trees that do not  
24 incorporate linear regression models at terminal nodes. Variable importance within the final  
25 model is measured by recording the increase in out-of-sample MSE that results when a  
26 covariate is randomly permuted for each tree in the ensemble. This increase in error is then  
27 averaged across all trees in the ensemble. In our model, the largest increases in error resulted  
28 from permutation of land cover and temperature, followed by 2-month lagged temperature  
29 and precipitation. Covariate influence can be evaluated through the use of partial dependence  
30 plots, which measure the change in model predictions that result from changing the value of  
31 one parameter while leaving all other covariates constant (Hastie et al., 2009). Partial  
32 dependence plots indicate that model predictions of streamflow are higher when the percent of

1 agricultural land cover is greater than approximately 75%, when temperatures anomalies are  
2 low, and when precipitation anomalies are high. However, it appears that the plot for lagged  
3 temperature might be sensitive to outliers at high temperature anomalies as evidenced by the  
4 large increase that occurs above an anomaly of +2, in a region where very few data points are  
5 present.

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

21 The partial dependence plots for the percentage of the basin classified as agricultural  
22 land cover indicates a positive relationship between agricultural land cover and streamflow in  
23 all basins except for the Gilgel Abbay (Fig. 7). This would be expected if deforestation had  
24 contributed to a decrease in evapotranspiration in the contributing watersheds. The exact  
25 nature of this response differs across the different rivers, with the relatively minor responses  
26 in Koga and Ribb, and much stronger responses in the Gumara and Megech basins. However,  
27 this plot also demonstrates some of the limitations associated with different model structures.  
28 The plot for Gumara is highly erratic, indicating that the M5 model might be overfit to the  
29 training dataset, despite the use of model averaging to reduce model variance. Additionally,  
30 the GAM used in the Megech basin was only trained on agricultural land cover values up to  
31 77%; while this model may be accurately representing the impact of land cover changes

1 within this range, extrapolating this relationship to higher values leads to predictions that may  
2 not be physically realistic.

### 3 **3.3 Climate Change Sensitivity and Uncertainty Assessment**

4 Fig. 8 shows the results of the climate change sensitivity analysis for total flow from all  
5 five tributaries, with dashed lines representing 95% confidence intervals obtained through 100  
6 bootstrapped resamples of the data set. As would be expected, increasing temperature  
7 independently of precipitation results in decreasing total flows while increasing precipitation  
8 results in higher flows. However, the uncertainty surrounding temperature sensitivity  
9 increases at higher changes in temperature, while the uncertainty surrounding precipitation  
10 sensitivity remains relatively constant, even at extreme changes in annual precipitation. The  
11 bottom panels of the figure show the sensitivity of total inflows to concurrent changes in  
12 temperature and precipitation. Unsurprisingly, decreasing precipitation combined with higher  
13 temperatures results in greater decreases in total flow than when temperature and precipitation  
14 are varied independently. However, even if temperature increases are combined with higher  
15 precipitation, total flows decline in the majority of bootstrap resamples.

16 The uncertainty surrounding temperature sensitivity is a key limitation to using data-  
17 driven approaches for climate impact assessment. To better understand which models and  
18 basins are contributing to this uncertainty, Fig. 9 shows how the coefficient of variation (the  
19 standard deviation of predictions from all bootstrap samples divided by the mean of these  
20 predictions) varies as a function of temperature change in each basin. From this figure, it is  
21 apparent that the Megech model is by far the largest contributor to model uncertainty;  
22 however, it is not clear whether this contribution is due to model structure (the GAM model  
23 used for the Megech River) or characteristics associated with the basin itself. To investigate  
24 how different model structures contributed to this uncertainty, the bootstrap resampling  
25 procedure was used to assess uncertainty in streamflow predictions in the Gumara River from  
26 all model types. This basin was chosen because all six models were able to outperform the  
27 climatology model, and thus could be considered good choices for model selection based on  
28 predictive accuracy alone. The results indicate that the increase in uncertainty is highest, and  
29 increases non-linearly, in the GLM, GAM, and MARS models. Uncertainty increases more  
30 slowly in the ANN and M5 models, and no noticeable increase in uncertainty is apparent in  
31 the random forest model.

## 1 **4 Discussion**

2           The objective of this study was not to identify the “best” approach for empirical  
3 rainfall-runoff modeling, as this is likely to be highly specific to the basin and problem to  
4 which a model is applied. However, we hope that the comparison conducted here can  
5 highlight some of the strengths and limitations of different approaches, as well as demonstrate  
6 some important issues that should be kept in mind for model comparisons in the future. One  
7 important finding was the limitation with using NSE as an error metric. Our results confirm  
8 previous studies that found that even uninformative models able to capture basic seasonality  
9 are able to achieve high NSE values (Legates and McCabe, 1999; Schaefli and Gupta, 2007),  
10 and provide further evidence indicating that high NSE values should be considered a  
11 necessary but not sufficient requirement for model usage in planning situations. For instance,  
12 the simple climatology model used for comparison purposes here is able to achieve high NSE  
13 values, but would be unsuitable for planning since it does not account for any interannual  
14 variability nor the possibility for non-stationary conditions caused by changing climate and  
15 land cover. In particular, understanding error structure can be valuable in evaluating whether  
16 model biases might undermine the model’s suitability for management activities. In our  
17 example, the autocorrelation present in the standard-formulation models meant that these  
18 models were consistently underestimating wet-season flows, resulting in low estimates of the  
19 total annual flow in the rivers. Since multiple reservoirs are planned for construction on these  
20 rivers to support irrigation activities, this bias could lead to poor estimates of how much water  
21 is available for agricultural use in the short term (ie., seasonal forecasting) and long-term (due  
22 to climate change). Interestingly, difficulties in accurately capturing high flows has been  
23 observed in physical hydrologic models for Ethiopia (e.g., Setegne et al., 2011; Mekonnen et  
24 al., 2009) and more generally (e.g., Wilby, 2005). The implications of this limitation should  
25 be carefully evaluated before using models for water resource planning or (more importantly)  
26 flood risk evaluation.

27           Depending on the model type used, different mechanisms are available to evaluate  
28 covariate importance and influence within the model. This evaluation can be useful in  
29 confirming that the model is replicating relationships between input and output variables in a  
30 reasonable manner. While the relationships identified in this evaluation are fairly  
31 straightforward (for example, increasing runoff with higher precipitation and lower  
32 temperatures), these simple relationships are still important in highlighting the mechanisms by

1 which the models make predictions so that they are not “black boxes.” For instance, Han et al.  
2 (2007) explore how ANN flood forecasting models responds to a double-unit input of rain,  
3 finding that some formulations respond in a hydrologically meaningful way to increased  
4 rainfall intensity, while others do not. Similarly, Galelli and Castelletti (2013a) describe how  
5 input variable importance can be used to highlight differences in hydrologic processes  
6 between an urbanized and forested watershed. The easy manner in which covariate  
7 relationships within the GAM and random forest models can be visualized using a single  
8 command within their respective R packages is a strong advantage to these approaches  
9 compared to methods such as M5 model trees and artificial neural networks. Of course, partial  
10 dependence plots can be developed for any model type (as was done in this research), but  
11 code must be written by the user and thus requires a higher degree of effort than is necessary  
12 for in-package functions. A downside to most machine-learning models is that they do not  
13 support the statistical formalism in assessing variable importance that is possible when linear  
14 models and GAMs are used. However, this formalism often rests on assumptions regarding  
15 model residuals that are unlikely to be met in many hydrologic models (Sorooshian and  
16 Dracup, 1980).

17         Within the Lake Tana basin, evaluation of covariate influence indicates that each  
18 basin’s model is performing in a reasonable manner, with runoff increasing with higher  
19 precipitation levels and decreasing with higher temperatures. The influence of precipitation  
20 and temperature is greatest in the current month, and progressively declines to a very small  
21 influence after two months. This suggests that long-term (multi-month) storage does not  
22 significantly contribute to variability in flow volumes. One interesting finding is the non-  
23 linear relationship between concurrent month precipitation and runoff that exists in the  
24 Megech and Ribb basins, which suggests that above a certain point increasing rainfall does  
25 not result in a commiserate increase in streamflow. Other studies have noted the dampening  
26 effect that wetlands and floodplains have had on river flows in the region (Dessie et al., 2014;  
27 Gebrehiwot et al., 2010); this phenomenon could explain the non-linear relationship identified  
28 in this work. The clearly negative relationship between temperature and runoff demonstrates  
29 the degree to which upstream evapotranspiration impacts streamflow and suggests that  
30 evapotranspiration is largely energy-limited, rather than water-limited. Increasing agricultural  
31 land-use appears to be associated with higher runoff in all rivers except for Gilgel Abbay  
32 (where no clear relationship between land cover and runoff was observed), and suggests that  
33 agricultural expansion at the expense of forest cover has reduced the evaporative component

1 of the water balance in these basins. Finally, the relative performance of different model  
2 formulations themselves can also be informative. For instance, the improved performance of  
3 the anomaly-formulation models indicates that the relationship between precipitation and  
4 runoff varies throughout the year and could point towards differences in runoff-generating  
5 mechanisms in the wet and dry seasons that have been observed in other case studies (Wilby,  
6 2005).

7         One limitation with data-driven approaches for streamflow prediction is that the  
8 relationships they model can only generate reliable predictions for conditions that are  
9 comparable to those experienced historically. Using these models to generate predictions for  
10 conditions that exceed historic variability is likely to introduce considerable uncertainty into  
11 their projections. Our results indicate that uncertainty in projections of streamflow under  
12 changing precipitation is relatively constant, whereas uncertainty increases markedly in  
13 projections of streamflow under increasing temperature. This result is not surprising when one  
14 considers the basin's climate, which is characterized by highly variable rainfall but fairly  
15 consistent temperatures (Table 6). A temperature increase of 3° C equates to almost two  
16 standard deviations beyond the historic mean, whereas a change in precipitation of 30% is  
17 well within the range of conditions experienced historically. One would expect that in other  
18 climates (for example, temperate watersheds with only minor changes in rainfall throughout  
19 the year), this relationship could be reversed. Despite the uncertainty that exists in projections  
20 of streamflow under changing temperature, total annual flow appears to be quite sensitive to  
21 increasing temperatures. In fact, the decreases in streamflow due to increasing temperature  
22 appears likely to be more than enough to counteract any increases in streamflow resulting  
23 from higher precipitation that is projected for the region in some global circulation models  
24 (GCMs). This is consistent with the work of Setegne et al. (2011), who used projections from  
25 multiple GCMs as input for a SWAT model developed for the region and found that  
26 streamflow decreased in the majority of emissions scenarios and models, even when  
27 precipitation increased. Unfortunately, this suggests that any hopes for a “windfall” of  
28 additional water to support agriculture and hydropower in the region under climate change  
29 may be unfounded.

30         Repeating the climate change sensitivity experiment with multiple models fit to the  
31 Gumara watershed indicated that the MARS, GAM, and linear models all result in the largest  
32 increase in uncertainty at high temperatures. This indicates that when models are fit to slightly

1 different bootstrap resamples of the historic dataset, the projected changes in streamflow at  
2 high temperature changes can be highly erratic. This is likely due to the fact that extrapolating  
3 the relationships that are observed between historic temperature and streamflow to higher  
4 temperatures can lead to very large changes in streamflow. Fitting the models to bootstrap  
5 resamples of the data results in minor changes to these relationships that can result in widely  
6 varying projections when the models are used to predict streamflow at higher temperatures,  
7 particularly when these relationships are nonlinear (as in the GAM). At the other end of the  
8 spectrum, the random forest model exhibits almost no increase in uncertainty at high  
9 temperatures, meaning that projections of streamflow at high temperatures are consistent  
10 across the bootstrap resamples. This is likely the result of the random forest model structure.  
11 The predicted value for each of a regression tree's terminal nodes is the average of all  
12 observations that meet the conditions described for that node. Thus, the model will not predict  
13 values beyond those experienced historically, even if covariate values exceed those contained  
14 within the historic dataset. Thus, this model is likely to underestimate the change in  
15 streamflow that results from increasing temperatures.

## 16 **5 Conclusions**

17 In this work, we compared multiple methods for data-driven rainfall-runoff modeling  
18 in their ability to simulate streamflow in five highly-seasonal watersheds in the Ethiopian  
19 highlands. Despite the popularity of ANNs in research on streamflow prediction to date,  
20 ANNs were not found to be the most accurate model in any of the five basins evaluated. Other  
21 methods, in particular GAMs and random forests, are able to capture non-linear relationships  
22 effectively and lend themselves to simpler visualization of model structure and covariate  
23 influence, making it easier to gain insights on physical watershed functions and confirm that  
24 the model is operating in a reasonable manner. However, it is important to carefully evaluate  
25 model structure and residuals, as these can contribute to biased estimates of water availability  
26 and uncertainty in estimating sensitivity to potential future changes in climate. In particular,  
27 autocorrelation in model residuals can result in underestimation of aggregate metrics such as  
28 annual flow volumes, even in models with high NSE performance. Uncertainty in GAM  
29 projections was found to rapidly increase at high temperatures, whereas random forest  
30 projections may be underestimating the impact of high temperatures on river flows. Thorough  
31 consideration of this uncertainty and bias is important any time that models are used for water  
32 planning and management, but especially crucial when using such models to generate insights

1 about future streamflow levels. By considering multiple model formulations and carefully  
2 assessing their predictive accuracy, error structure and uncertainties, these methods can  
3 provide an empirical assessment of watershed behavior and generate useful insights for water  
4 management and planning. This makes them a valuable complement to physical models,  
5 particularly in data-scarce regions with little data available for model parameterization, and  
6 warrants additional research into their development and application.

## 7 **Acknowledgements**

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

20

21



## 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&cpsid=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 Schaefli, 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 <sup>2</sup> )	Average annual streamflow at gauge (MCM)	Standard deviation of annual streamflow (MCM)	Coefficient of variation of annual streamflow	Average temp (°C)	Average monthly rainfall [mm]	
						May-Oct	Nov-Apr
Gilgel Abbay	2664	1883	217	0.12	15.7	206	39.3
Gumara	385	236	71	0.30	17.7	186	29
Koga	200	114	31	0.27	15.7	206	39.3
Megech	424	172	66	0.31	20.6	234	41.4
Ribb	677	210	83	0.36	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	family = Gaussian	NA
GAM	mgcv	family = Gaussian method = generalized cross validation variable selection = true basis dimension $k = 3$ epsilon = $10^{-7}$ maxit = 200	
MARS	earth	nk = 21 thresh = 0.001 fast.k = 20 pmethod = backward	degree = {1, 2, 3} nprune = {5, 10, 15, 20, 25}
ANN	nnet	weights = 1 rang = 0.7 maxit = 100 maxNWts = 1000 abstol = $10^{-4}$ reltol = $10^{-8}$	size = {1, 2, 4, 8, 20} decay = {0.0, 0.1, 0.5, 1.0, 2.0}
RF	randomForest	ntree = 500 sampsize = 528 nodesize = 5 nPerm = 1	mtry = {2, 3, 4, 5, 6, 7}
M5	Cubist	rules = 100 extrapolation = 100 sample = 0	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 highest performing  
3 standard formulation and anomaly formulation models in each basin (with model type in  
4 parenthesis), and resulting mean annual observed and predicted flow.

5

	Autocorrelation Factors		Mean Annual Flow (MCM)		
	Standard	Anomaly	Observed	Standard	Anomaly
Gilgel	0.33 (RF)	0.11 (GLM)	22,925	20,703	22,958
Gumara	0.29 (RF)	0.07 (M5)	2,870	2,392	2,734
Koga	0.04 (M5)	0.10 (GLM)	1,383	1,333	1,386
Megech	0.05 (RF)	0.04 (GAM)	2,035	1,637	2,028
Ribb	0.21 (RF)	-0.01 (RF)	2,575	1,969	2,615

6

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.

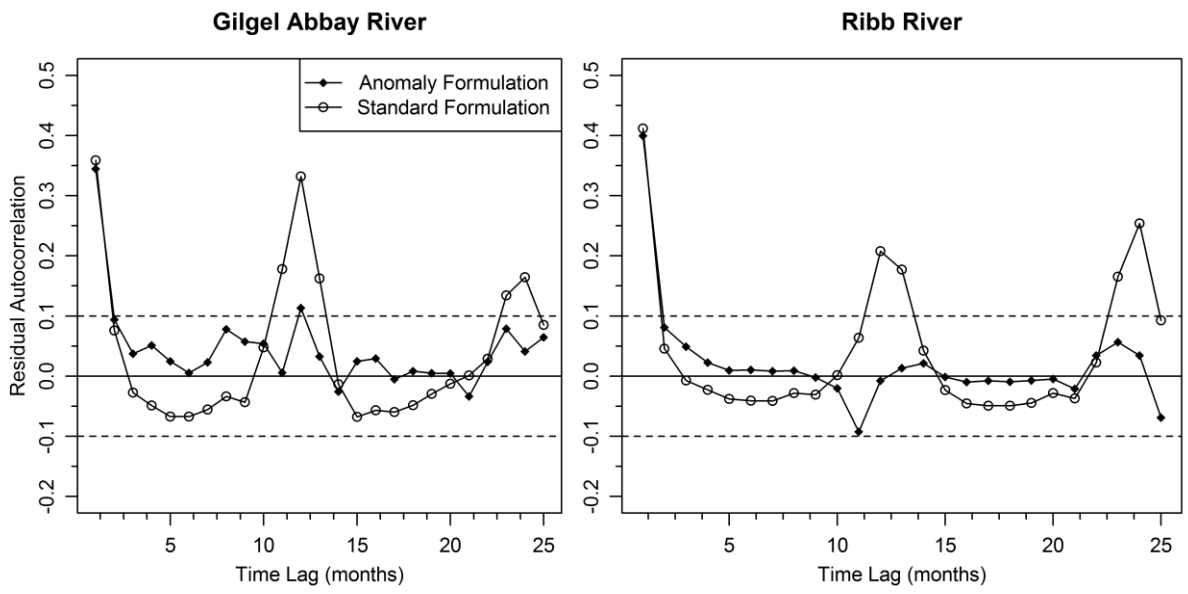
3

	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



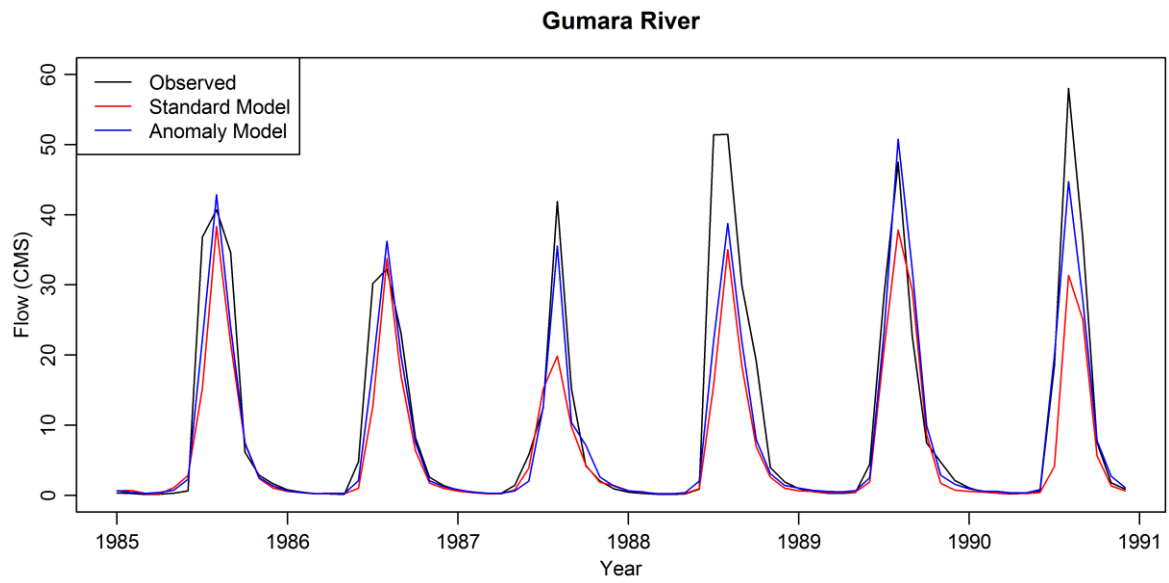
1 Figure 2. Autocorrelation in model residuals for the Gilgel Abbay and Ribb Rivers



2

3

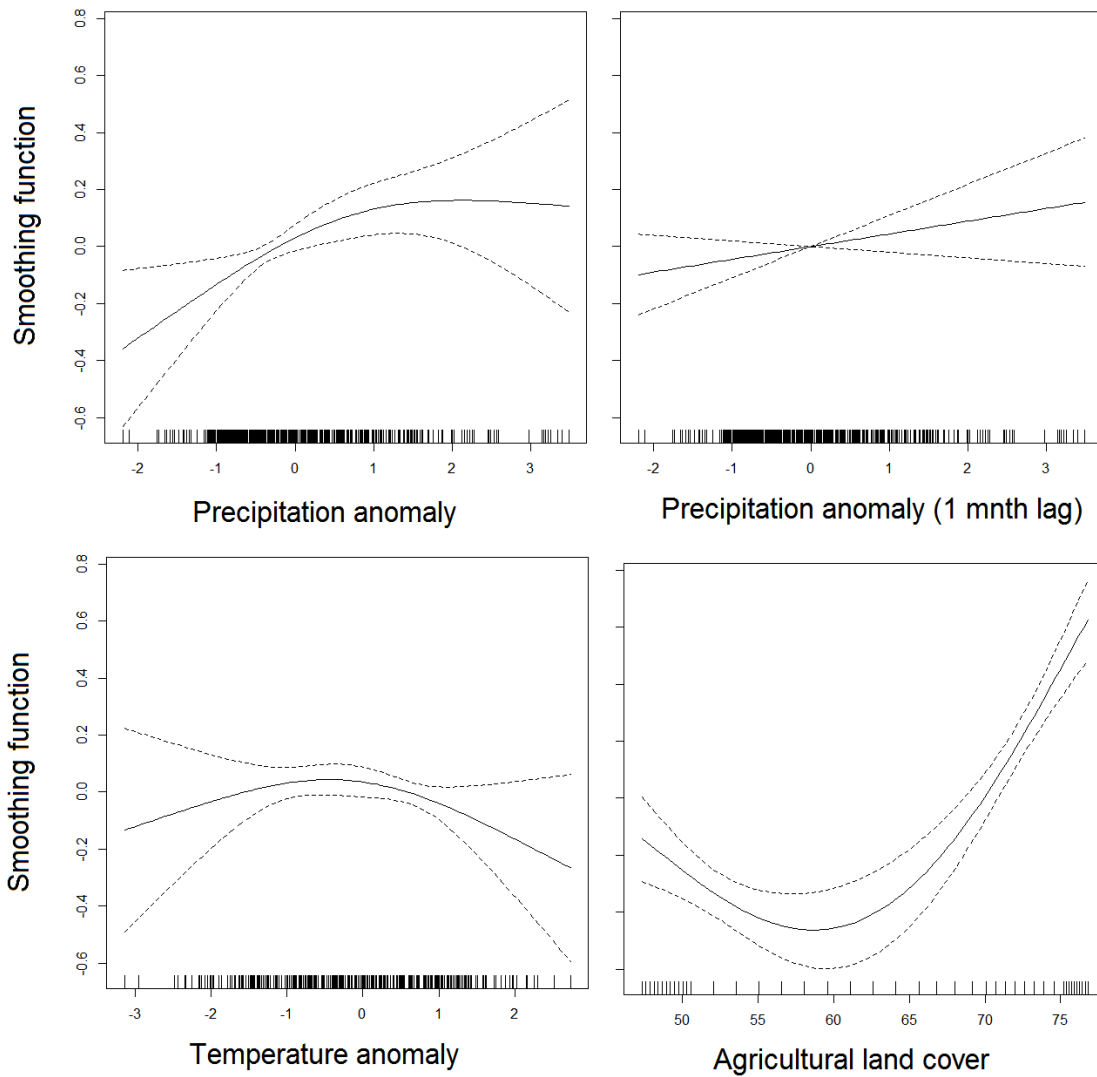
1 Figure 3. Example observed and predicted flows from the standard formulation RF model and  
2 anomaly formulation M5 model for the Gumara River from 1985 to 1991.



3

4

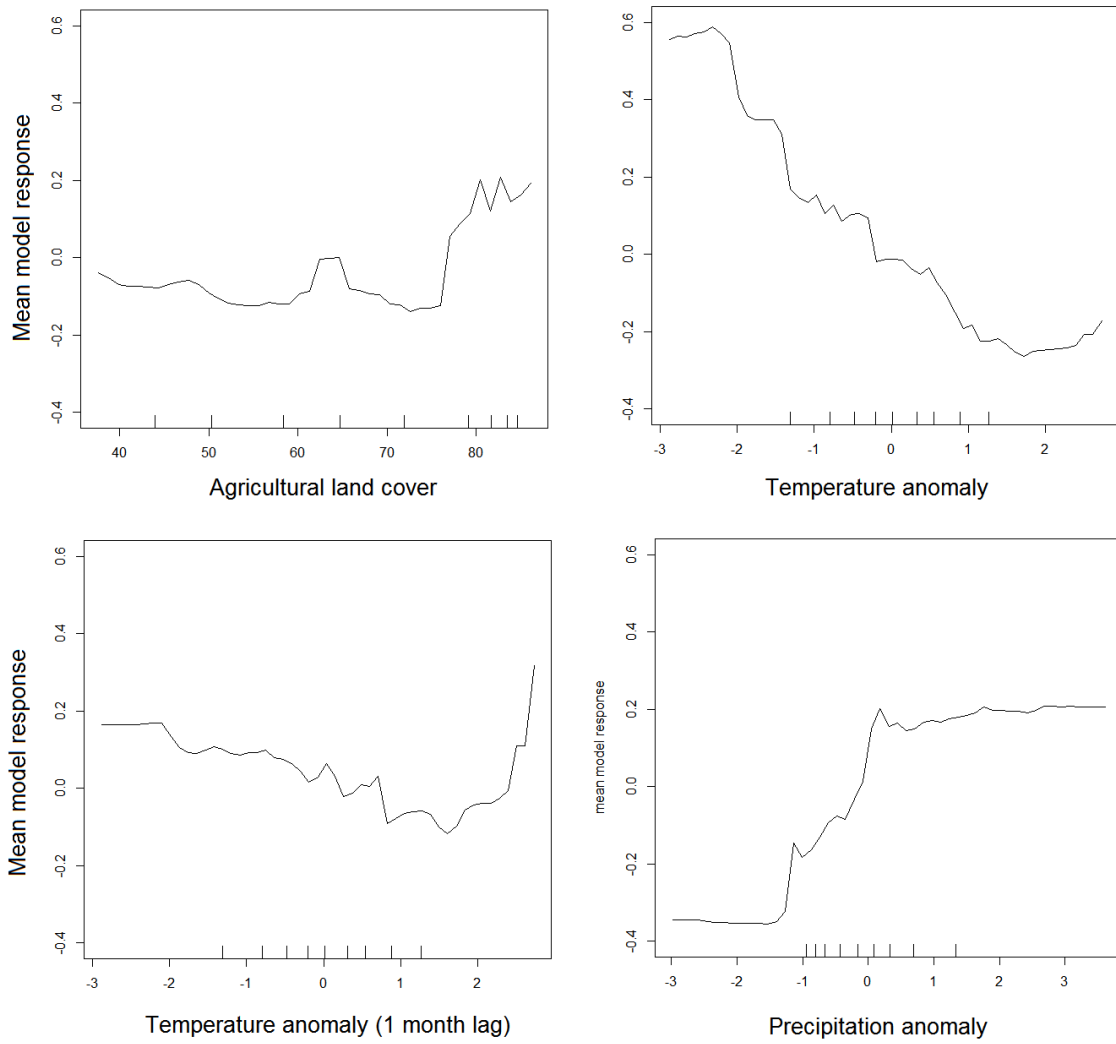
1 Figure 4. Plots of the smoothing functions used in the Megech River GAM. Hash marks along  
2 the x-axis indicate observation values of each covariate.



3  
4

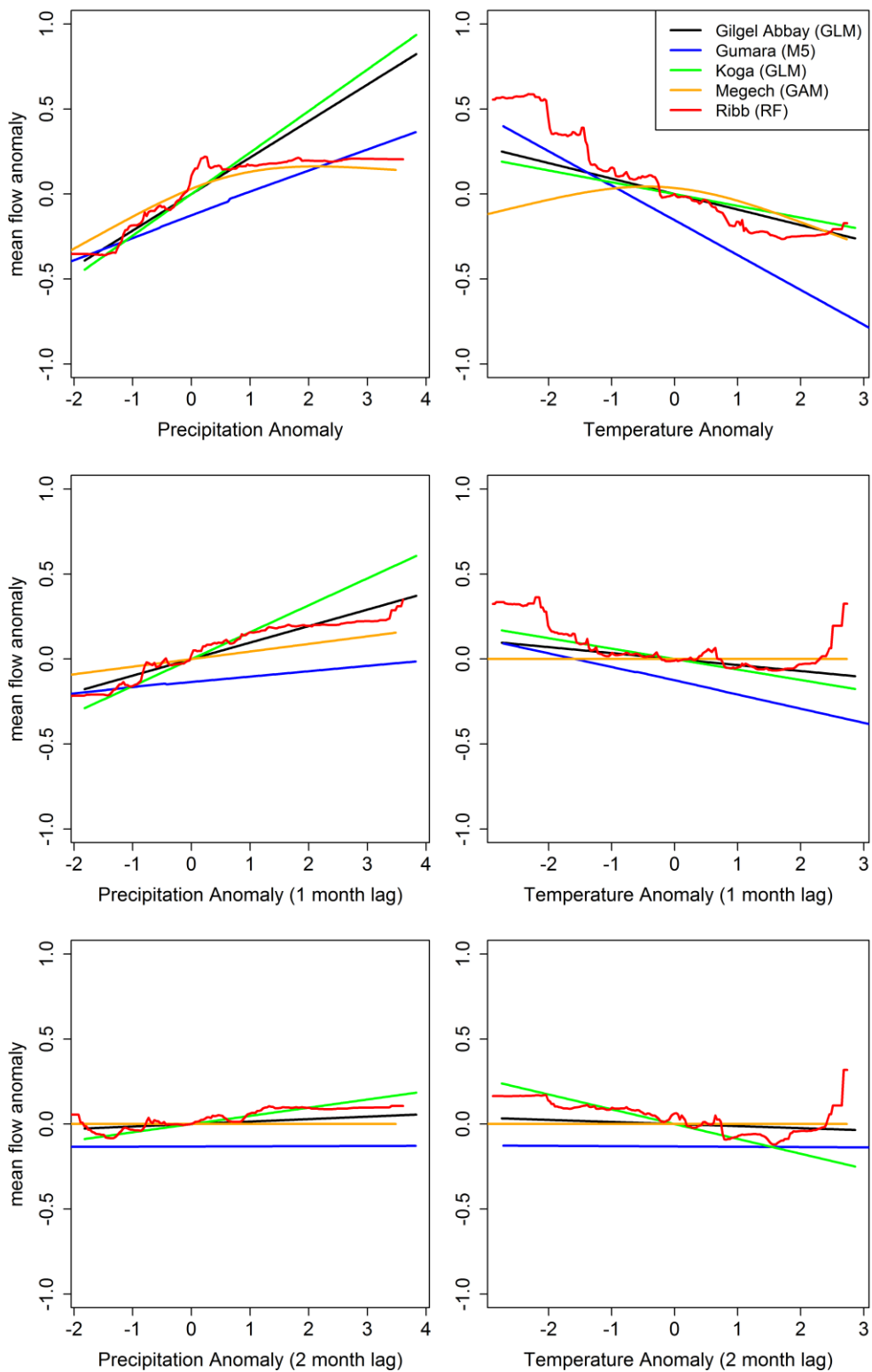


1 Figure 5. Partial dependence plots for the Ribb River random forest model. Hash marks along  
2 the x-axis show covariate sample decile values.



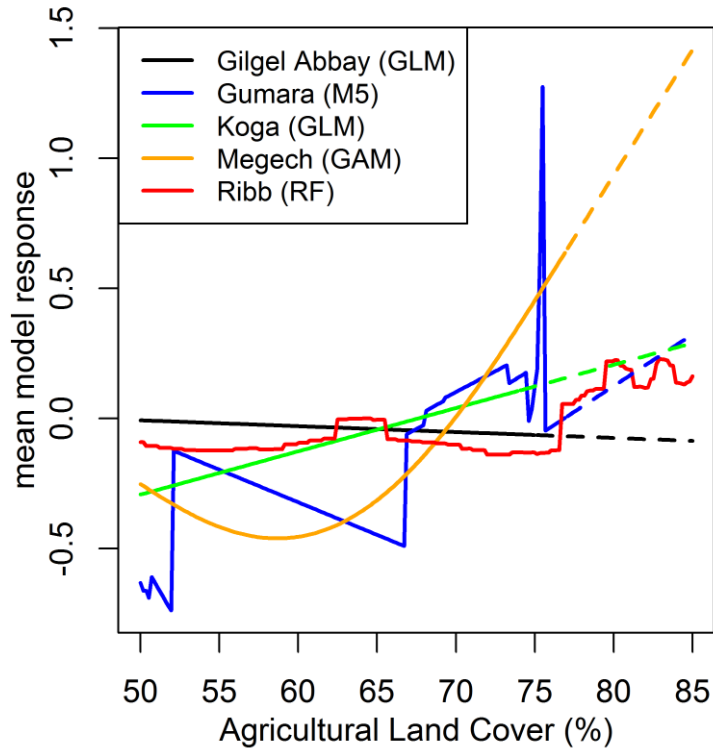
3  
4

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



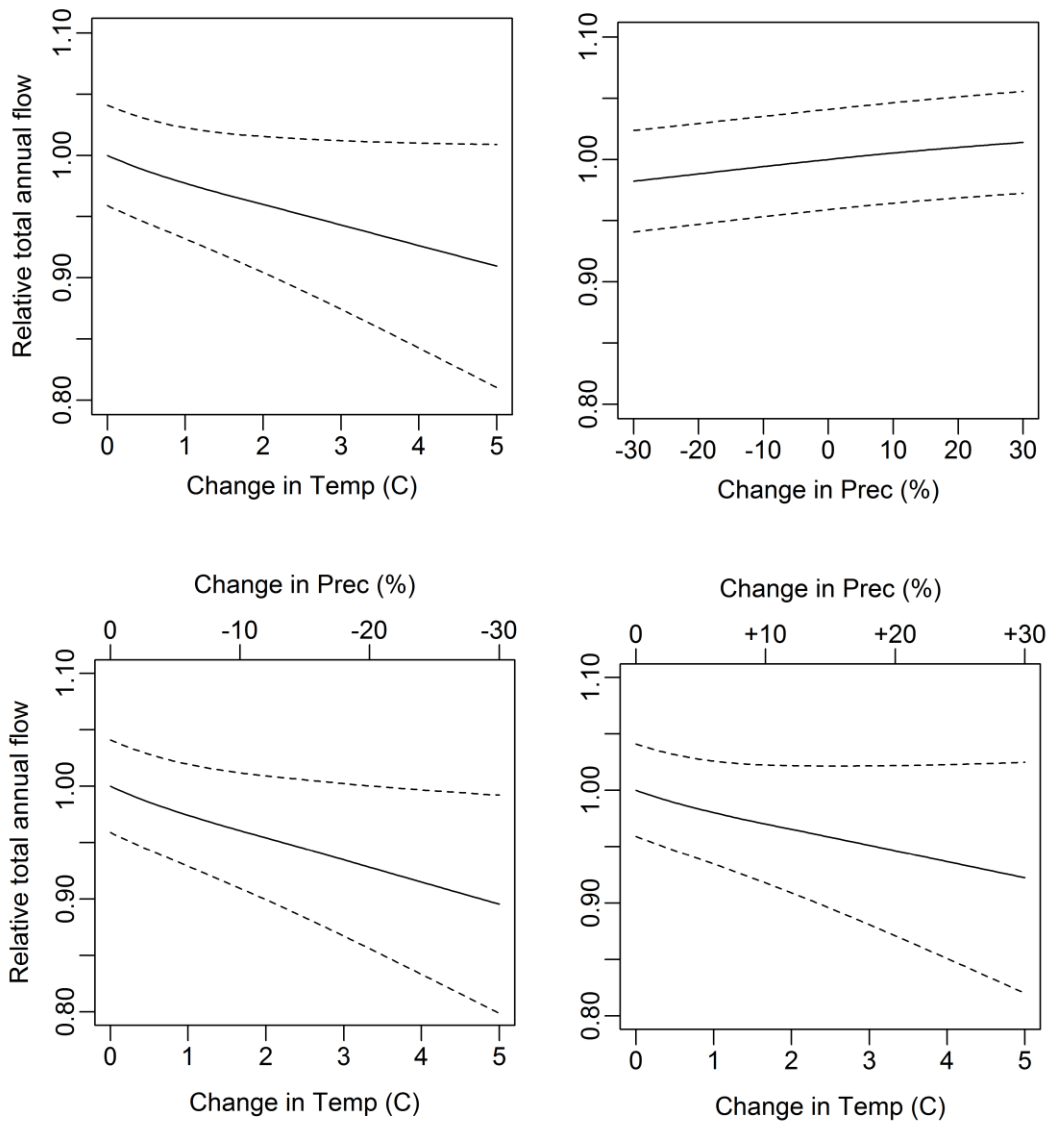
3

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.



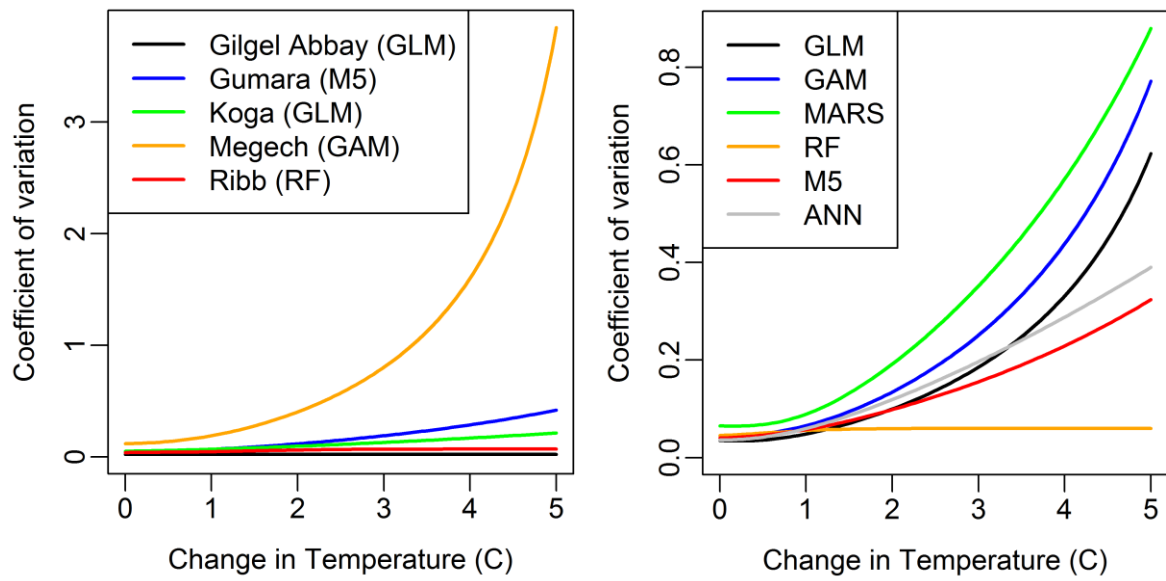
5  
6

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