1	A calibration method for projecting future
2	extremes via a linear mapping of parameters
3	Jeongjin Lee ^{1*} , Daniel Cooley ² , Anna M. Wagner ³ , Glen E. Liston ⁴
4	^{1*} Department of Mathematics and Statistics, Lancaster University,
5	Fylde College, Lancaster, LA1 4YF, United Kingdom.
6	² Department of Statistics, Colorado State University, Fort Collins,
7	80523. Colorado, United States of America.
8	³ U.S. Army Engineer Research and Development Center, Cold Regions
9	Research and Engineering Laboratory. Fort Wainwright, Alaska, United
10	States of America
11	⁴ Cooperative Institute for Besearch in the Atmosphere Colorado State
12	University, Fort Collins, Colorado, United States of America.
13 14 15	*Corresponding author(s). E-mail(s): j.lee58@lancaster.ac.uk; Contributing authors: cooleyd@stat.colostate.edu; Anna.M.Wagner@erdc.dren.mil; Glen.Liston@colostate.edu;
16	Abstract
17	In order to study potential impacts arising from climate change, future projec-
18	tions of numerical model output often must be calibrated to be comparable to
19	observations. Rather than calibrating the data values themselves, we propose a
20	novel statistical calibration method for extremes that assumes there exists a linear
21	relationship between parameters associated with model output and parameters
22	in both parameter estimates and the linear calibration, which we achieve via
23	bootstrap. To focus on extreme behavior, we assume both model output and
25	observations have distributions composed of a mixture model combining a Weibull
26	distribution with a generalized Pareto distribution for the tail. A simulation study
27	shows good coverage rates. We apply the method to project future daily-averaged
28	river runoff at the Purgatoire River in southeastern Colorado.

29 Keywords: climate projections, downscaling, flooding, extremes, calibration

³⁰ 1 Introduction

Numerical models are widely used in the Earth sciences to study atmospheric and 31 ocean dynamics, hydrology, atmospheric chemistry, and other processes. Because they 32 are driven by the known physics of the studied system, numerical models are the 33 best tools available to produce projections under possible climate scenarios. However, 34 numerical model output can have notable discrepancies from observations, which can 35 have implications for quantifying potential climate change impacts. These discrepan-36 cies (or 'model bias' in the climate literature) can occur because the model's spatial 37 support (often a grid cell) differs from that of point-referenced observations and/or 38 because the model's physics are necessarily a simplification and cannot capture the 39 full complexity of the Earth system. To better understand the behavior of quanti-40 ties of interest in relation to observations, there is a need to calibrate (or downscale) 41 model output. Various statistical calibration and bias correction methods have been 42 developed to address these issues. We provide an overview of widely applied methods 43 in Section 2.2. As part of their review, Teutschbein and Seibert (2012) compare dif-44 ferent correction methods specifically for regional climate model (RCM) simulations 45 in hydrological impact studies. 46

The calibration problem can be visualized as in Figure 1. The boxes with green checkmarks have data (either model output or observations) available. The red 'X' indicates there are no observations under the projected climate. Quantities of interest related to projected observations must be estimated based on a modeled relationship between model output and observations learned from the historical period, which is then applied to the projected period.



Fig. 1: Illustration of the calibration method applied to projected observations under the projected climate. Boxes with green checkmarks have data available and quantities of interest can be estimated directly. The red 'X' indicates there are no observations under the projected climate, and quantities of interest must be estimated based on the relationship between model output and observations, and the relationship between historical and projected climate.

⁵³ Our particular calibration study is motivated by a project which aims to estimate ⁵⁴ potential flood risk from the Purgatoire River to infrastructure at a military base

in southeastern Colorado, USA. For this river, flows are often greatest in the spring 55 due to runoff from melting snow, and it is of interest to know how climate-induced 56 changes to the timing and duration of snow cover could affect flood risk. Notably, 57 discrepancies between model output and local observations are amplified when focus-58 ing on extremes, as shown in Figure 2. With the growing attention on climate change 59 impacts, numerous studies have contributed to the development of calibration meth-60 ods for extremes to better assess their effects on local extremes (e.g., Schubert and 61 Henderson-Sellers, 1997; Vrac and Naveau, 2007; Benestad, 2010). In a comparison 62 of advanced downscaling methods, Bürger et al. (2012) evaluated approaches for 63 extremes, including automated regression-based statistical downscaling (ASD) (Hes-64 sami et al., 2008), bias correction spatial disaggregation (BCSD) (Wood et al., 2002), 65 and quantile regression neural networks (QRNN) (Taylor, 2000). In the context of 66 machine learning techniques, Campozano et al. (2016) compared statistical downscal-67 ing methods with two machine learning methods, specifically artificial neural network 68 (ANN) and least squares support vector machines (LA-SVM), to evaluate downscaled 69 general circulation model (GCM) estimates of monthly precipitation. 70

In this study, given the limited and relatively short data records, and without the 71 inclusion of additional predictor variables, our calibration method is classified as a 72 transfer function approach, as opposed to stochastic weather generators and weather 73 typing methods (Vrac and Naveau, 2007). Accordingly, in Section 4.3, we focus on 74 comparing calibration methods that emphasize the direct relationships between large-75 scale model output and local observations. We propose a novel univariate calibration 76 method for extremes, with the development of a multivariate version left for future 77 work. 78

Our primary aim is to provide estimates of high quantiles of projected Purgatoire 79 River observations. In particular, we wish to provide estimates (with uncertainty) of 80 quantiles roughly corresponding to the 1-in-10 and 1-in-100 year events, the latter 81 of which will require extrapolation into the tail as the data records we employ are 82 much shorter. The need for extrapolation leads us to employ a parametric model. 83 Our model, which will be fit to the entire distribution, will rely on an extreme value 84 model to capture the behavior in the upper tail. Unlike the advanced downscaling 85 methods referenced in Bürger et al. (2012), where data outside the range of the fitted 86 quantiles are extrapolated using specific distributions such as Weibull or exponential 87 distributions, we fit a generalized Pareto distribution to the tail, allowing for more 88 flexibility in capturing different tail behaviors. 89

Our calibration method assumes there exists a linear transfer function governing the relationship between the parameters of the distributions describing model output and the observations. The historical period will be used to estimate this linear relationship, and then it is applied to the parameters of the projected model output to obtain an estimate of the distribution of projected observations. In addition to accounting for the uncertainty associated with parameter estimates, a bootstrap method will allow us to also account for uncertainty associated with estimating the transfer function.

In contrast to dynamical downscaling using RCMs, statistical calibration requires
thorough validation, especially for extremes, to verify the methods as noted by (Bürger
et al., 2012). We perform validation using the historical period, given that observations

for projected climates are not available. To evaluate the performance of calibration
 methods for extremes, we use quantile-quantile plots (QQ-plots), a standard tool in
 extreme value analysis, and summary statistics derived from large values exceeding a
 high threshold.

Section 2 describes the data and reviews common statistical calibration and bias correction methods. In Section 3.1, we motivate and describe our model for both observed and modeled river flows. Section 3.2 describes the process of our calibration method. In Section 3.3, we describe how we obtain uncertainty estimates for the model parameters, transfer function, and estimated quantiles. We then present simulation results, case study, and method comparisons in Section 4. Finally, we conclude with a summary and discussion.

2 Data and other calibration methods

112 2.1 Data description

Throughout the study, we analyze daily-averaged river flow observations and river flow 113 model output for the Purgatoire River in Colorado covering the period from 2002 to 114 2013. Projected river flow model output is also produced for this period under a pro-115 jected climate (not shown in Figure 2). Complete measurements are taken across three 116 datasets, resulting in a sample size of n = 1,836. The superimposed timeseries plot 117 of the observations and model output for the historical period is shown in Figure 2a. 118 Flows are only shown for the period between April and August as this is the period 119 when the river is at risk for flooding, and river flow measurements outside these months 120 are considered unreliable. 121



Fig. 2: (a) Timeseries plot of daily-averaged river discharges (black solid line) and modeled river discharges (blue dashed line) for the Purgatoire River in Colorado for 2002-2013. (b) QQ-plot comparing empirical quantiles of observations to empirical quantiles of model output for the same historical period.

The process to produce modeled river flows is quite involved. Two simulations 122 were forced with a high-resolution (4km grid) dataset obtained from the Weather 123 Research and Forecasting (WRF) model. The current climate simulation was forced 124 with ERA-Interim reanalysis data (Dee et al., 2011) for the period from October 125 2000 to September 2013. The projected pseudo-global warming simulation (PGW) (a 126 perturbation experiment) for the same period is forced with ERA-Interim reanalysis 127 and a climate perturbation (Rasmussen et al., 2011). This perturbation reflects the 128 95-year mean change signal from the Coupled Model Intercomparison Project Phase 129 5 (CMIP5) multi-model ensemble under the high-emissions RCP8.5 scenario. The 130 ensemble-mean monthly climate change was derived from 19 CMIP5 models selected 131 based on their performance, as detailed in Section 2.2. of Liu et al. (2017). 132

This high-resolution modeled weather output was used as input for SnowModel (Liston and Elder, 2006; Liston et al., 2020), a numerical model for the accumulation, evaporation, and melting of snow over a study area. Finally, modeled snow runoff was combined with WRF-produced meteorology and input into HydroFlow (Liston and Mernild, 2012), a numerical hydrological model which produces simulated streamflow measurements.

The modeled stream runoff is useful for understanding changes in timing and relative runoff amounts between the historical and projected periods, but there is a clear mismatch between the distributions of the modeled river flows and the observations. In particular, the extremes of the observations are not represented by the modeled runoff. Calibration is needed to use the modeled river flow output for assessment of projected flood risk.

¹⁴⁵ 2.1.1 Model evaluation

To assess the goodness of fit for numerical models such as SnowModel, several summary statistics can be used. We consider the determinant coefficient, $R^2 \in [0,1]$, where values close to 1 indicate a better fit; root mean squared error, RMSE ≥ 0 , with lower values implying a better fit; and Nash-Sutcliffe efficiency coefficient (Nash and Sutcliffe, 1970), NSE $\in (-\infty, 1]$, which reflects the proportion of the variance in the observations that is accounted for by the model relative to the total variance of the observations

NSE = 1 -
$$\frac{\sum_{t=1}^{n} (x_{Mod}(t) - x_{Obs}(t))^2}{\sum_{i=1}^{n} (x_{Obs}(t) - \bar{x}_{Obs})^2}$$
,

where $x_{Mod}(t)$ and $x_{Obs}(t)$ are the model output and observations at time $t = 1, \ldots, n$, 153 respectively, and \bar{x}_{Obs} denotes the sample mean of the observations. An NSE of 1 indi-154 cates a perfect match between model output and observations. While these measures 155 are typically useful for evaluating mean behavior, our focus is on high-quantiles. To 156 make these statistics more relevant for extremes, we also consider large values exceed-157 ing a high threshold. Setting the high threshold at the 0.95 quantile for each dataset 158 of observations and model output for the historical period, the obtained statistics are 159 summarized in Table 1. As a graphical diagnostic, we create a QQ-plot of observa-160 tions versus modeled river flows in Figure 2b, showing a clear mismatch in the higher 161 quantiles. 162

	\mathbb{R}^2	RMSE	NSE
SnowModel	0.14 / 0.10	4.64 / 13.85	-0.014 / -0.524

Table 1: Summary statistics for SnowModel under the historical climate. Values before the slash are derived from all data points, while values after the slash are derived from data exceeding a high threshold.

¹⁶³ 2.2 Statistical calibration and bias corrections

We review commonly applied statistical calibration and bias correction methods. Let 164 $\boldsymbol{x}_{Mod}^{h}(t)$ denote the variables from model output for the historical period at time t, 165 let $x_{Qbs}^{h}(t)$ be the observations of variable of interest from the historical period, and 166 let $x_{Mod}^p(t)$ represent the model output for the projected period. Statistical calibra-167 tion methods take these available data to attempt to describe $x_{Obs}^{p}(t)$, the unavailable 168 observations for the projected period. The model output used does not have to rep-169 resent the projected variable of interest; $\boldsymbol{x}_{Mod}^{p}(t)$ might be larger-scale predictors, 170 for example principal components of geopotential heights (e.g., Hanssen-Bauer et al., 171 2005). Huang et al. (2019) say that most proposed statistical calibration methods fall 172 into three general categories: regression, a shift/scale (or delta method) approach, or 173 quantile mapping. We note that the choice of bias correction method often depends 174 on the specific motivations, data characteristics, and underlying assumptions of the 175 study. 176

177 2.2.1 Regression-based method

¹⁷⁸ A generic form of regression is assumed as

$$x_{Obs}^{h}(t) = f(\boldsymbol{x}_{Mod}^{h}(t), \boldsymbol{\beta}) + \epsilon(t).$$

Often the function f is standard linear regression. Calibration via regression makes the most sense when the observations and model output are synchronous; that is, the modeled weather at time t represents the actual weather at t. Climate reanalysis data is synchronous, but output from general circulation models is typically not. When the data is asynchronous, then the distribution of the model output needs to be related to the distribution of the observations. A notable study that uses the regression-based downscaling method includes (Wilby et al., 1999).

¹⁸⁶ 2.2.2 A shift/scale method (delta method)

A simple shift/scale approach is moment based. Letting $X^h(t)$ and $X^p(t)$ be the random variables representing the variable of interest in the historical period and projected period respectively and \cdot can be either model output or observations, the shift/scale approach assumes

$$X^p_{\cdot}(t) = s(X^h_{\cdot}(t) + m).$$

¹⁸⁷ Thus, the projected distribution function is shifted by m and scaled by s:

$$F_{X_{\cdot}^{p}(t)}(x) = F_{X_{\cdot}^{h}(t)}((x-m)/s)$$

The parameters m and s are learned from the historical period and applied to the projected period. As an example, a univariate version of linear bias correction for climate model output by Bürger et al. (2012) is

$$\hat{x}_{Obs}^{p}(t) = \left(\frac{x_{Mod}^{p}(t) - \bar{x}_{Mod}^{p}}{\sigma_{Mod}^{h}}\right) \sigma_{Obs}^{h} + (\bar{x}_{Mod}^{p} - \bar{x}_{Mod}^{h}) + \bar{x}_{Obs}^{h},$$
(1)

¹⁹¹ where \bar{x}_{Mod}^{p} , \bar{x}_{Mod}^{h} , and \bar{x}_{Obs}^{h} represent the sample means of the model output for the ¹⁹² project period, the model output for the historical period, and the observations for ¹⁹³ the historical period, respectively, and σ_{Obs}^{h} and σ_{Mod}^{h} are the standard deviations of ¹⁹⁴ observations and model output for the historical period. A study by Teutschbein and ¹⁹⁵ Seibert (2012) compares various shift and scale corrections, including linear scaling, ¹⁹⁶ variance scaling, power transformation, and the delta-change method.

This shift/scale method is simple and m and s are easily estimated, but is best suited for understanding changes in the center or bulk of the distribution and is of limited value for detecting changes in extremes.

²⁰⁰ 2.2.3 Quantile mapping (QM)

Quantile mapping (QM) defines the transfer function that connects cumulative distribution functions (CDFs) of F_{Obs}^h and F_{Mod}^h

$$\hat{x}_{Obs}^{p}(t) = \hat{F}_{Obs}^{h^{-1}} \{ \hat{F}_{Mod}^{h}[x_{Mod}^{p}(t)] \},$$
(2)

where $\hat{x}_{Obs}^{p}(t)$ corresponds to the projected observations at time t within the projected period, \hat{F}_{Obs}^{h} and \hat{F}_{Mod}^{h} are estimated CDFs of observations and model output under historical climate. As the quantile mapping transfer function uses only the information from the historical period, it can fail for extreme values if $x_{Mod}^{p}(t)$ falls outside the range of values used to estimate \hat{F}_{Mod}^{h} . Hence, extrapolation is required for extremes.

²⁰⁸ 2.2.4 Quantile delta mapping (QDM)

Instead of using the direct extrapolation, there have been methods such as equidistant and equiratio quantile mapping (Li et al., 2010; Wang and Chen, 2014), which use the information from the CDF of the projected model output, denoted by F_{Mod}^{p} , for the projected period. Cannon et al. (2015) verified that those quantile mapping approaches are equivalent to the quantile mapping of the 'delta change method' (Olsson et al., 2009) and termed this approach 'quantile delta mapping' (QDM).

QDM preserves relative changes (or deltas) in all quantiles between the projected model output $x_{Mod}^{p}(t)$ and historical model output $x_{Mod}^{h}(t)$. The first step involves detrending the projected model output by quantile, followed by bias correction of the projected model output to historical observations via quantile mapping. After that,

the relative changes, denoted by $\Delta_{Mod,Rel}(t)$, in the quantiles of the projected model output are applied to the bias-corrected historical values to capture the projected climate change signal (or relative changes in modeled quantiles)

$$\hat{x}_{Obs}^{p}(t) = \hat{F}_{Obs}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \} \times \Delta_{Mod,Rel}(t)$$

$$= \hat{F}_{Obs}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \} \times \frac{x_{Mod}^{p}(t)}{\hat{F}_{Mod}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \}}$$
(3)

To preserve absolute changes, denoted by $\Delta_{Mod,Abs}(t)$, in quantiles (e.g., temperatures in Celsius, see Cannon et al., 2015), additive deltas are applied to historical bias-

²²⁴ corrected values. Specifically, the adjusted quantile is given by $\widehat{F}_{Obs}^{h^{-1}}\{\widehat{F}_{Mod}^{p}[x_{Mod}^{p}(t)]\}$

$$\hat{x}_{Obs}^{p}(t) = \hat{F}_{Obs}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \} + \Delta_{Mod,Abs}(t) = \hat{F}_{Obs}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \} + x_{Mod}^{p}(t) - \hat{F}_{Mod}^{h^{-1}} \{ \hat{F}_{Mod}^{p}[x_{Mod}^{p}(t)] \}$$
(4)

²²⁵ For more details, refer to Cannon et al. (2015).

Compared to the methods described above, our calibration method in Section 3.2 employs a parametric approach for extrapolation. We compare these methods and evaluate their performance in Section 4.3.2.

²²⁹ 3 Methodology

3.1 A mixture model for modeling the distribution's bulk and tail

We create a model for the the entire distribution of projected observations, and rely on extreme value models to characterize the upper tail. To motivate our eventual model, we briefly describe modeling and calibrating only the extremes.

An aim of an extreme value analysis is to not let data in the bulk of the distri-235 bution influence estimates of tail behavior. Thus, classical extremes methods analyze 236 only an extreme subset of the data: either block (e.g., annual) maxima or threshold 237 exceedances. Here, our data record is insufficient to model only annual maxima, and 238 seasonal effects would make fitting sub-annual maxima (e.g., monthly) dubious. Con-239 sider a standard peaks-over-threshold approach, in which the generalized Pareto (GP) 240 distribution (Balkema and De Haan, 1974) is fit to data which exceed a high threshold 241 u. The GP distribution is defined as 242

$$G(x; u, \xi, \sigma) = P(X < x \mid X > u) = \begin{cases} 1 - \left[1 + \frac{\xi}{\sigma}(x - u)\right]_{+}^{-1/\xi}, & \xi \neq 0, \\ 1 - \exp\left[-\left(\frac{x - u}{\sigma}\right)\right]_{+}, & \xi = 0, \end{cases}$$
(5)

where $z^+ = \max(z, 0)$. The GP density is

$$g(x; u, \xi, \sigma) = \sigma^{-1} \left\{ 1 + \frac{\xi}{\sigma} \left(x - u \right) \right\}^{-1/\xi - 1},$$
(6)

on the support $\{x : x > u, 1 + \frac{\xi}{\sigma}(x-u) > 0\}$ with a scale parameter $\sigma > 0$ and a shape parameter $\xi \in (-\infty, \infty)$. The shape parameter ξ determines the fundamental nature of the tail. If $\xi < 0$, the distribution of threshold excesses has a bounded tail with an upper end point $u < x < u - \sigma/\xi$. The zero shape parameter $\xi = 0$ corresponds to a light tail (where g becomes the exponential distribution in the limit), and the positive shape parameter $\xi > 0$ corresponds to a heavy tail.

The threshold is selected (not estimated), typically using diagnostic plots. In order to perform calibration, the threshold for the observations and model output would need to have a known relationship, and one approach would be to employ thresholds defined by a common exceedance probability.

Although our primary interest is in estimating very high quantiles, we desire to perform calibration of the entire distribution so that any quantity of interest can be projected. We consider a model which flexibly fits the distribution's tail, which behaves like a generalized Pareto in the upper tail, and which restricts data in the bulk from influencing parameter estimates in the tail. Our mixture model assumes that for a threshold u > 0, the probability density function of the random variable X is

$$h(x;\boldsymbol{\theta}) = d(\boldsymbol{\theta})^{-1} \Big[(1 - \pi(x;u,\delta)) \frac{f(x;\beta,\lambda)}{F(u;\beta,\lambda)} (1 - \kappa_u) + \pi(x;u,\delta) g(x;u,\xi,\sigma) \kappa_u \Big], \quad (7)$$

 $(\beta, \lambda, \xi, \sigma)$ is the parameter vector, $f(x;\beta,\lambda)$ where θ = 260 $(\beta/\lambda)(x/\lambda)^{\beta-1}\exp(-(x/\lambda)^{\beta})$ for x > 0 represents the Weibull density with scale 261 parameter $\lambda > 0$ and shape parameter $\beta > 0$, $F(u; \beta, \lambda)$ is the Weibull CDF evaluated 262 at $u, \kappa_u = P(X > u)$, and $g(x; u, \xi, \sigma)$ is the generalized Pareto density in (6). The 263 weight function 264

$$\pi(x; u, \delta) = \begin{cases} 0, & \text{if } x < u \\ \frac{1}{\delta}, & \text{if } u \le x \le u + \delta \\ 1, & \text{if } x > u + \delta \end{cases}$$
(8)

for $\delta > 0$ provides a continuous transition between the models for the bulk and tail. The normalizing constant $d(\boldsymbol{\theta}) = \int_0^\infty h(x; \boldsymbol{\theta}) dx$ is needed only because of the weight function; if $\delta = 0$, then $\pi(x; u, 0)$ is defined as 0 for x < u and 1 otherwise, implying $d(\boldsymbol{\theta}) = 1$ for any parameter values.

Given independent and identically distributed (i.i.d.) observations x_1, \ldots, x_n , infer-269 ence begins by choosing an appropriate threshold u and transition range parameter 270 δ . The threshold u can be chosen via visual diagnostics such as the mean resid-271 ual life plot (Davison and Smith, 1990) and δ can be chosen qualitatively so that 272 the transition between the bulk and the tail is satisfactory. With u and δ selected, 273 $\hat{\kappa}_u = n^{-1} \sum_{i=1}^n \mathbb{1}(x_i > u)$, and numerical maximum likelihood (ML) can then be 274 performed to find estimates for β, λ, ξ , and σ , where $d(\theta)$ is calculated by numerical 275 integration. 276

There have been models proposed which fit the entire distribution, but which also 277 have an upper tail which behaves asymptotically like a generalized Pareto. Often 278 the aim of these models is to avoid the issues associated with selecting a threshold. 279 One method, proposed by Papastathopoulos and Tawn (2013) and extended by both 280 Naveau et al. (2016) and Stein (2021) constructs a model via a composition of the 281 generalized Pareto distribution $G(x; u, \xi, \sigma)$ in (5) and another 'carrier' CDF, Q(v)282 for $v \in [0, 1]$. Specifically, the key motivation in the Naveau et al. (2016) model is as 283 follows: we can simulate from the GP distribution by applying a uniformly distributed 284 random sample U into the GP quantile function $G^{-1}(U)$. By replacing U with a more 285 flexible random variable $V = Q^{-1}(U)$, where Q is a continuous CDF on [0, 1], the 286 resulting random variable $Y = G^{-1}[Q^{-1}(U)]$ forms a more flexible distribution family. 28 They propose a class of CDFs, Q(v), such that the upper tail retains the behavior of 288 the GP distribution, and the CDF of Y near zero behaves like a power function y^{κ} . 289 Four parametric families satisfying these conditions are introduced, and for simplicity, 290 we focus on the first two for model comparison. 291

²⁹² 1. $Q(v) = v^{\kappa}, \quad \kappa > 0, v \in [0, 1]$ ²⁹³ 2. $Q(v) = pv^{\kappa_1} + (1 - p)v^{\kappa_2}, \quad \kappa_1, \kappa_2 > 0, v, p \in [0, 1].$

Another approach, termed a 'mixture' by Scarrott and MacDonald (2012), com-294 bines a density model for the bulk with a generalized Pareto density for the tail, often 295 smoothing the transition between them. A particular mixture model was proposed by 296 Frigessi et al. (2002) which used a Weibull model for the bulk and a GPD for the tail. 297 In contrast to our weight function (8), their model uses a Cauchy CDF as a weight 298 function to transition between the bulk and the tail over the entire data range. The 299 GP shape parameter in the Frigessi et al. (2002) model tends to be either overesti-300 mated or underestimated, as noted by Naveau et al. (2016), and the scale parameter 301 of the Cauchy CDF, which controls the transition speed, is challenging to estimate. 302

We investigate the tail behavior of these parametric models without using threshold selection and assess the goodness of fit in Section 4.3, comparing it to our fixed threshold approach.

306 3.2 Calibration method via linear mapping of distribution 307 parameters

Our calibration method assumes the linear relationship between parameters associated with model output and parameters associated with observations. Specifically, we will assume

$$\boldsymbol{\theta}_{Obs}^{\cdot} = A\boldsymbol{\theta}_{Mod}^{\cdot} + \boldsymbol{b},\tag{9}$$

where θ_{Obs} and θ_{Mod} denote the parameter vectors of the distributions for the observations and the model output respectively, and ' \cdot ' is a placeholder for both h denoting the historical and p denoting the projected climate. The historical period will be used to estimate this linear relationship, and then it is applied to the parameters of the projected model output to obtain an estimate of the distribution of projected observations.

³¹⁷ With the available observations and model output, one can obtain ML estimates for ³¹⁸ $\boldsymbol{\theta}_{Obs}^{h} = (\beta_{Obs}^{h}, \lambda_{Obs}^{h}, \xi_{Obs}^{h}, \sigma_{Obs}^{h}), \boldsymbol{\theta}_{Mod}^{h}$, and $\boldsymbol{\theta}_{Mod}^{p}$, the parameter vectors for the extreme

³¹⁹ mixture model applied to observations in historical climate, model output under his-³²⁰ torical climate, and model output under projected climate, respectively. Furthermore, ³²¹ asymptotic estimates for the respective covariance matrices Σ_{Obs}^{h} , Σ_{Mod}^{h} , and Σ_{Mod}^{p} ³²² can be obtained using the inverse of the observed Fisher information matrix. The cal-³²³ ibration method uses these estimates to produce an estimate of $\boldsymbol{\theta}_{Obs}^{p}$, the parameter ³²⁴ vector for projected observations, for which there are no observations to analyze.

Let A be a 4×4 matrix and **b** be a 4×1 vector. Both A and **b** are assumed time-invariant; that is, this same linear relationship holds both under historical and projected climate. Essentially, this assumption says that the model 'biases' which lead to the discrepancy between the modeled river flow output and the observations do not depend on the climate state; without some similar assumption, the calibration problem is impossible.

One can use parameter estimates for both the observations and model output in the historical climate to estimate A and b. Based on (9), we consider the second-order moment form, motivated by the asymptotic normality of ML estimators

$$\Sigma^h_{Obs} = A \Sigma^h_{Mod} A^\top.$$
⁽¹⁰⁾

Solving (10) with estimates plugged in yields $\hat{A} = \hat{\Sigma}_{Obs}^{h^{1/2}} \hat{\Sigma}_{Mod}^{h^{-1/2}}$ and plugging into (9) yields $\hat{b} = \hat{\theta}_{Obs}^{h} - \hat{\Sigma}_{Obs}^{h^{1/2}} \hat{\Sigma}_{Mod}^{h^{-1/2}} \hat{\theta}_{Mod}^{h}$. We obtain square root matrices by spectral decomposition so that generically $\Sigma^{1/2} = P \Lambda^{1/2} P^{\top}$, where P is the square orthogonal matrix of Σ 's eigenvectors and Λ is the diagonal matrix of eigenvalues. If the covariance matrix is not positive definite due to closely corrected ML estimates or numerical issues, it is necessary to obtain the nearest positive definite covariance matrix before performing the spectral decomposition (e.g., using the nearPD function from the Matrix R package).

With estimates \widehat{A} and \widehat{b} , the parameter estimate for the projected observations is

$$\hat{\boldsymbol{\theta}}_{Obs}^{p} = \widehat{A}\hat{\boldsymbol{\theta}}_{Mod}^{p} + \hat{\boldsymbol{b}}.$$
(11)

343 3.3 Uncertainty quantification for projected parameter a44 estimates

The quantity of primary interest is $\hat{\theta}^p_{Obs}$ and we devise a bootstrap procedure to quan-345 tify its uncertainty which arises from the uncertainty of parameter estimators as well 346 as that of the estimated linear projection in (11). Let $x_{Obs}^{h}(t), x_{Mod}^{h}(t), x_{Mod}^{p}(t), t =$ 347 $1, \ldots, n$, be observations in historical climate, model output under historical climate, 348 and model output under projected climate, respectively. For each dataset, resample with replacement to obtain bootstrap ML estimates, $\hat{\theta}^{h\,(b)}_{Obs}$, $\hat{\theta}^{h\,(b)}_{Mod}$, $\hat{\theta}^{p\,(b)}_{Mod}$, $\hat{\Sigma}^{h\,(b)}_{Obs}$, $\hat{\Sigma}^{h\,(b)}_{Mod}$, $\hat{\Sigma}^{h\,(b)}_{Obs}$, $\hat{\Sigma}^{h\,(b)}_{Mod}$, 349 350 and $\hat{\Sigma}_{Mod}^{p(b)}$, where $b = 1, \dots, B$, the number of bootstrap iterations. One can obtain $\hat{A}^{(b)}$ and $\hat{b}^{(b)}$ as in Section 3.2 to obtain $\hat{\theta}_{Obs}^{p(b)}$. Bootstrap confidence intervals can then 351 352 be constructed. 353

To account for serial dependence in real data applications such as the river flow data for the Purgatoire River, we apply the block bootstrap method. The approach involves spliting the data into non-overlapping blocks and resamples these blocks to preserve the temporal dependence in each block. To determine the block length, we
use both the blockboot function from the OBL R package, which selects the optimal
block length in terms of the minimum root mean squared error, and the autocovariance
function (ACF).

361 4 Results

³⁶² 4.1 Simulation study

We assess coverage rates via a simulation study. Values for $\boldsymbol{\theta}_{Obs}^{h}$, $\boldsymbol{\theta}_{Mod}^{h}$, $\boldsymbol{\theta}_{Mod}^{p}$ are set to the ML estimates for our application (given in the first three rows of Table 3 below). Assuming (9) and using covariance estimates from the real data, we solve for A, b, and obtain $\boldsymbol{\theta}_{Obs}^{p}$.

For each simulation iteration, n = 5,000 i.i.d. realizations are drawn by acceptreject algorithm from the extreme mixture models for historical observations, historical model output, and projected model output. For each sample, we draw B = 1,000bootstrap samples and obtain bootstrap ML estimates for these three distributions. $\hat{A}^{(b)}, \hat{b}^{(b)}$, and $\hat{\theta}^{p(b)}_{Obs}$ are obtained as described in Section 3.2.

Bootstrap based 95% confidence intervals are produced for all parameter estimates. Parallel computing is used to repeat the simulation 100 times, and Table 2 reports coverage rates. Coverage rates for the quantity of interest θ_{Obs}^{p} appear reasonable, especially given the coverage rates for the other parameters which do not require calibration.

	Coverage Rates
$ \begin{array}{c} (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{h}_{Mod} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{b}_{Obs} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{p}_{Mod} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{p}_{Mod} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{p}_{Ob} \end{array} $	(0.97, 0.98, 0.97, 0.95) (0.92, 0.93, 0.93, 0.95) (0.97, 0.94, 0.85, 0.93) (0.99, 0.99, 0.95, 0.97)

Table 2: Coverage rates for 95% boot-strap based confidence intervals for 100simulations.

4.2 Case study: application to river discharges

We apply the calibration method to daily-averaged runoff model output and river discharges from the Purgatoire River in southeastern Colorado. We use complete measurements from three data sets: historical model output, projected model output, and historical river discharges for the period from April to August between 2002 and 2013, as described in Section 2.1.

To assess the tail behavior of the historical observations, we first fit a GP distribution to the observations exceeding the 0.95 quantile ($q_{0.95} = 10.27$), obtaining the shape parameter estimate of $\hat{\xi}^h_{Obs} = 0.05$. The 95% confidence interval for ξ is (-0.21, 0.32), suggesting the possibility of either a light tail or a slightly heavy tail.

We then fit a mixture model $h(x; \theta)$ in (7) to each data set, obtaining the ML estimates reported in the first three rows of Table 3. The threshold u is set to $q_{0.95}$,

	ML estimates	Standard Errors
$ \begin{array}{c} (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{h}_{Mod} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{h}_{Obs} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{p}_{Mod} \\ (\hat{\beta}, \hat{\lambda}, \hat{\xi}, \hat{\sigma})^{p}_{Obs} \end{array} $	$\begin{array}{c}(1.31,2.76,-0.13,2.18)\\(0.57,0.91,0.05,8.99)\\(1.27,2.42,0.01,3.10)\\(\textbf{0.53},\textbf{0.57},\textbf{0.19},\textbf{9.92})\end{array}$	$\begin{array}{c}(0.03,0.07,0.10,0.39)\\(0.01,0.05,0.11,1.78)\\(0.03,0.05,0.12,0.58)\\(\textbf{0.03},\textbf{0.11},\textbf{0.45},\textbf{6.30})\end{array}$

Table 3: ML estimates for parameters and their correspond-ing bootstrap-based standard errors.

388

the empirical 0.95 quantile of each data set and $\delta = q_{0.96} - u$. Figure 3 shows QQ-plots for the mixture model fit to river discharges in historical climate, model output under historical climate, and model output under projected climate. Of particular interest is the upper tail, and the QQ-plots show a reasonable fit accounting for the usual model uncertainty associated with estimating extreme behavior.

The estimate $\hat{\xi}^{h}_{Obs} = 0.05$ agrees with our preliminary generalized Pareto fit show-394 ing that the bulk data do not influence the tail estimate. The estimate of $\hat{\xi}^h_{Mod} = -0.13$ 395 from the historical model output, with a 95% confidence interval of (-0.37, 0.07), may 396 suggest a bounded tail. This is because the confidence interval contains more nega-397 tive plausible values when accounting for sampling uncertainty. The bounded tail for 398 model output may illustrate the challenges that the sequence of numerical models face 399 in replicating extreme behavior. The GP shape parameter estimate from the projected 400 model output is $\hat{\xi}_{Mod}^p = 0.01$, and eight instances of modeled river flows under the 401 projected climate exceed the maximum modeled river flow under the historical climate 402 as shown in the center and right panels of Figure 3. To test whether the tail behav-403 ior of the model output under the projected climate is significantly greater than that 404 under historical climate, a one-sided Wald test was conducted. The resulting p-value of 405 0.21 indicates that the data do not provide sufficient evidence to support a significant 406 difference in the tail behavior of the model output between the two climates. 407

We next find the linear relationship between river discharges and model output 408 in historical climate through \hat{A} and \hat{b} and then apply the linear relationship under 409 projected climate to obtain $\hat{\theta}_{Obs}^p$, reported in the fourth row of Table 3. Differences 410 between $\hat{\theta}_{Obs}^{p}$ and $\hat{\theta}_{Obs}^{h}$ reflect the differences between $\hat{\theta}_{Mod}^{p}$ and $\hat{\theta}_{Mod}^{h}$, as both the GP shape and scale parameters have increased. The estimate $\hat{\xi}_{Obs}^{p} = 0.19$ seems large, 411 412 but its standard error of 0.45 reflects the uncertainty in estimating tail indices with 413 short data records, and the uncertainty in projecting observations. Figure 4 shows the 414 densities associated with the fitted models for both model output and observations 415 under both historical and projected climate, and shows kernel density estimates for 416 the three combinations with data. 417

We report 95% block-bootstrap confidence intervals for parameters of θ_{Obs}^{h} and θ_{Obs}^{p} , respectively in the first four columns of Table 4. Using the blockboot tool and



Fig. 3: (a) QQ-plot of river discharges under historical climate; (b) QQ-plot of model output under historical climate; and (c) QQ-plot of model output under projected climate.

ACF together, we set the block lengths to 48, 30, and 30 days for historical obser-420 vations, historical model output, and projected model output, respectively. The 95% 421 block-bootstrap confidence intervals for (ξ, σ) in a GP distribution are of primary 422 interest. Not surprisingly, we observe wider bootstrap confidence intervals for $(\xi, \sigma)^p$ 423 than ones for $(\xi, \sigma)^h$ due to the additional uncertainty of linear projection. Estimates 424 of high quantiles are of more practical interest than the parameter estimates. We pro-425 vide 95% block-bootstrap confidence intervals for the 0.9993 and 0.99993 quantiles 426 of the historical and projected observations. These correspond to the 1-in-1500 and 427 1-in-15000 day event; as there are 152 observations in the April to August period, 428 these are roughly 1-in-10 and 1-in-100 year events. The width of the confidence inter-429 vals under the historical period is wide due to the short data record we employ, and 430 this uncertainty becomes amplified when projected. Nevertheless, despite the limited 431 information in the data and model runs for this risk study, it seems there is potential 432 for higher river flows and thus increased flood risk under the projected climate. 433

	β	λ	ξ	σ	$q_{0.9993}$	$q_{0.99993}$
$egin{array}{c} m{ heta}_{Obs}^h \ m{ heta}_{Obs}^p \end{array}$	(0.55, 0.60)	(0.83, 1.00)	(-0.13, 0.29)	(5.64, 12.31)	(40.36, 60.37)	(57.27, 1768.95)
	(0.51, 0.60)	(0.53, 0.85)	(-0.70, 0.90)	(6.70, 29.73)	(44.33, 344.51)	(49.37, 4357.51)

Table 4: 95% block-bootstrap based confidence intervals for parameters associated with local observations as well as for the 0.9993 and 0.99993 quantiles under historical climate and projected climate, respectively.

14



Fig. 4: Dashed line shows the fitted mixture model density (a) for model output under historical climate, (b) for observed river discharges under historical climate, (c) for model output under projected climate, and (d) for river discharges under projected climate. Kernel density estimates are shown with solid lines in (a), (b), and (c).

4.3 Method validation and comparative assessment with other approaches

4.3.1 Goodness of fit for other parametric models without threshold selection

We explore the tail behavior of other parametric extremes models described in Section 3.1, starting by fitting each model to the entire historical river discharges. We used the fit.extgp tool in the mev R package to fit the Naveau et al. (2016) model considering two carrier functions: $Q(v) = v^{\kappa}$ and $Q(v) = pv^{\kappa_1} + (1-p)v^{\kappa_2}, \kappa_1, \kappa_2 > 0, v, p \in$ ⁴⁴² [0, 1]. We also fit the Frigessi et al. (2002) model by numerical maximum likelihood ⁴⁴³ estimation. Both models produced heavy tail ML estimates with $\hat{\xi}^{h}_{Obs} = 0.56$ for the ⁴⁴⁴ Frigessi et al. (2002) model and $\hat{\xi}^{h}_{Obs} = 0.99$ for both carrier functions in the Naveau ⁴⁴⁵ et al. (2016) model.

For goodness of fit, QQ-plots of the empirical quantiles against the fitted model quantiles for the two carrier functions showed a clear mismatch between the modeled upper tail and the largest observations in Figure 5b and 5c. While the QQ-plot for the Frigessi et al. (2002) model in Figure 5a performed better, it still showed more discrepancies in the higher quantiles compared to our model with the fixed threshold approach in Figure 3a. A similar issue arises for the historical model output (not shown).



Fig. 5: QQ-plot of empirical quantiles of historical observations versus fitted model quantiles (a) for the Frigessi et al. (2002) model, (b) for the Naveau et al. (2016) model with a carrier $Q(v) = v^{\kappa}, v \in [0, 1]$, and (c) for the Naveau et al. (2016) model with a carrier $Q(v) = pv^{\kappa_1} + (1-p)v^{\kappa_2}, \kappa_1, \kappa_2 > 0, v, p \in [0, 1]$.

While the asymptotic tail behavior of these models follows the GP distribution in the limit, this case study indicates either the estimation is challenging, or that both models do not sufficiently separate the tail behavior from the bulk, resulting in data from the bulk unduly influencing the shape parameter estimate. Therefore, we opt to use a fixed threshold approach for the tail to ensure the proper calibration of extremes.

458 4.3.2 Validation for different calibration methods

Focusing on extremes, we perform validation for different statistical calibration meth-459 ods outlined in Section 2.2, using only the historical period, as projected observations 460 are unavailable. Ideally, both historical and projected datasets with sufficiently long 461 records would be used for a more reliable evaluation. However, in the absence of such 462 data, we split both the historical observations and model output into a calibration set 463 $(60\%, \text{ a sample size } n_{cal} = 1, 101)$ and a validation set $(40\%, \text{ a sample size } n_{val} = 735)$, 464 assuming that these datasets still preserve distinct climate characteristics. The cal-465 ibration set is used as the historical period to estimate parameters for the mixture 466

⁴⁶⁷ model in (7) and the transfer functions in Section 2.2, and the validation set serves as
⁴⁶⁸ the out-of-sample data to assess goodness of fit and evaluate the performance of the
⁴⁶⁹ calibration methods.

We compare projected observations (bias-corrected values) derived from different statistical calibration methods to actual observations (assumed to be observed) in the validation set. We consider six calibration methods: linear parameter mapping, QDM preserving absolute changes in quantiles (4), QDM preserving relative changes in quantiles (3), QM (2), the linear scaling approach (1), and a simple regression approach.

To properly evaluate these calibration methods for extremes, we use both QQplots and summary statistics that take datasets where both observations and projected observations exceed the 0.95 quantile, respectively. Similarly in Section 2.1.1, we report the summary statistics in Table 5. The values before the slash represent summary statistics calculated using all data points, while the values after the slash are derived from data exceeding the high threshold.

	\mathbb{R}^2	RMSE	NSE
Linear mapping	0.97 / 0.98	0.86 / 3.58	-1e-3 / -0.02
QDM-abs	0.38 / 0.70	4.20 / 10.85	-0.02 / 0.12
QDM-rel	0.39 / 0.71	3.79 / 9.68	-0.03 / -0.01
QM	0.38 / 0.68	3.38 / 10.55	-2e-3 / -0.14
Linear scaling	0.31 / 0.67	4.32 / 10.17	-0.03 / -0.01
Simple regression	0.29 / 0.68	3.57 / 13.15	-0.04 / -0.81

Table 5: Summary statistics of \mathbb{R}^2 , RMSE, and NSE for six different calibration methods. QDM-abs indicates QDM preserving absolute changes in quantiles, and QDM-rel stands for QDM preserving relative changes in quantiles.

For the linear mapping of parameters approach, we obtain GP shape parameter estimates of $\hat{\xi}^{h}_{Obs} = 0.09$ for the historical climate in the training set and $\hat{\xi}^{p}_{Obs} =$ 0.12 for the projected climate in the validation set. As the projected observations in the validation set are quantile-matched to the corresponding actual observations for the projected period, the correlation coefficient, or \mathbb{R}^2 , between projected and actual observations is effectively close to 1.

To visualize the performance of the six approaches for extremes, QQ-plots are 488 shown in Figure 6. The high quantile values from the linear mapping approach align 489 closely with the diagonal, indicating a satisfactory fit for extreme values. The QDM 490 results also appear reasonable, indicating potential for further improvement through 491 parametric estimation of the CDFs. It is important to note that this method com-492 parison does not definitely conclude that any one method is superior in all cases. 493 The performance of each method may vary depending on the characteristics of the 494 real data and the available predictor variables, where different transformations and 495 parameterizations could lead to different results. 496



Fig. 6: (a) QQ-plot of empirical quantiles of observations versus fitted model quantiles in the validation set for a linear mapping of parameters. QQ-plots of empirical quantiles of actual versus projected observations in the validation set (b) for a QDM preserving absolute changes in quantiles, (c) for a QDM preserving relative changes in quantiles, (d) for QM, (e) for a linear scaling approach, and (f) for a simple regression approach, respectively.

497 5 Summary and Discussion

We develop a novel statistical calibration method focusing on extreme values by 498 applying a linear parameter mapping approach rather than directly calibrating the 499 model output. This method has a fundamental assumption of a linear relationship 500 between parameters associated with model output and those with local observations. 501 Once estimated under historical climate conditions, this linear relationship is applied 502 to parameter estimates under projected climates. This stationarity assumption is 503 standard to most statistical calibration methods. However, given that both model 504 output and observations exhibit temporal dependence, the dependence could affect 505 uncertainty estimates. Given a long data record, a possible extension to this work is 506 to incorporate temporal dependence in the calibration method by considering non-507 508 stationarity in this linear relationship, that is, incorporating temporal dependence or possible covariates to the scale or shape parameters of the GP distribution. 509

18

To flexibly fit the tail of the distribution, we employ a mixture distribution with 510 an extreme value model for the tail; however, the proposed calibration method itself 511 does not require any particular distribution. We opt for a fixed threshold approach for 512 tail approximation due to a better performance for extremes compared to automatic 513 threshold selection methods. But, this fixed threshold approach introduces some sub-514 jectivity. Some other approach that can reduce the subjectivity of a threshold selection 515 would be useful for improving the calibration method for extremes. Our method is 516 currently applied in a univariate case, focusing on changes in the marginal distribu-517 tion. Extending this approach to a multivariate calibration using a multivariate GP 518 distribution could be an interesting future work. 519

Importantly, to account for uncertainty in both parameter estimates and linear 520 projection, a bootstrap approach is employed. Accounting for uncertainty is at least 521 as important as the projected quantity estimates themselves. Of course, the projec-522 tions we produce are based on only one possible climate scenario, and the uncertainty 523 of human behavior and its affect on the climate likely outweighs the estimation uncer-524 tainty we capture here. As the climate changes, there is increased need to produce 525 projections to inform planning for potential future outcomes. Numerical models are 526 powerful, but imperfect tools, and calibration methods such as ours can help to account 527 for when model output does not accurately reflect the quantities needed by decision 528 makers. 529

Acknowledgments JL, DC, AW, and GL's work was supported by Strategic
 Environmental Research and Development Program (SERDP) Resource Conservation
 (RC) and Resilience Project RC-2515. JL and DC were also partially supported by
 National Science Foundation (NSF) DMS-1811657.

534 **References**

- Balkema, A.A., De Haan, L.: Residual life time at great age. The Annals of probability,
 792–804 (1974)
- Benestad, R.E.: Downscaling precipitation extremes. Theoretical and Applied Climatology 100(1-2), 1-21 (2010)
- Bürger, G., Murdock, T., Werner, A., Sobie, S., Cannon, A.: Downscaling
 extremes—an intercomparison of multiple statistical methods for present climate.
 Journal of Climate 25(12), 4366–4388 (2012)
- Cannon, A.J., Sobie, S.R., Murdock, T.Q.: Bias correction of gcm precipitation by
 quantile mapping: How well do methods preserve changes in quantiles and extremes?
 Journal of Climate 28(17), 6938–6959 (2015)
- Campozano, L., Tenelanda, D., Sanchez, E., Samaniego, E., Feyen, J.: Comparison of
 statistical downscaling methods for monthly total precipitation: case study for the
 paute river basin in southern ecuador. Advances in Meteorology 2016(1), 6526341
 (2016)
 - 19

- Davison, A.C., Smith, R.L.: Models for exceedances over high thresholds. Journal of
 the Royal Statistical Society: Series B (Methodological) 52(3), 393-425 (1990)
- 551 Dee, D.P., Uppala, S.M., Simmons, A.J., Berrisford, P., Poli, P., Kobayashi, S., Andrae,
- U., Balmaseda, M., Balsamo, G., Bauer, d.P., et al.: The era-interim reanalysis:
- 553 Configuration and performance of the data assimilation system. Quarterly Journal

- Frigessi, A., Haug, O., Rue, H.: A dynamic mixture model for unsupervised tail estimation without threshold selection. Extremes 5(3), 219–235 (2002)
- Hanssen-Bauer, I., Achberger, C., Benestad, R., Chen, D., Førland, E.: Statistical
 downscaling of climate scenarios over scandinavia. Climate Research 29(3), 255–268
 (2005)
- Hessami, M., Gachon, P., Ouarda, T.B., St-Hilaire, A.: Automated regression-based
 statistical downscaling tool. Environmental modelling & software 23(6), 813–834
 (2008)
- Huang, Y.-N., Reich, B.J., Fuentes, M., Sankarasubramanian, A.: Complete spatial
 model calibration. The annals of applied statistics 13(2), 746–766 (2019)
- Liston, G.E., Elder, K.: A distributed snow-evolution modeling system (snowmodel).
 Journal of Hydrometeorology 7(6), 1259–1276 (2006)
- Liu, C., Ikeda, K., Rasmussen, R., Barlage, M., Newman, A.J., Prein, A.F., Chen,
 F., Chen, L., Clark, M., Dai, A., *et al.*: Continental-scale convection-permitting
 modeling of the current and future climate of north america. Climate Dynamics 49,
 71–95 (2017)
- Liston, G.E., Itkin, P., Stroeve, J., Tschudi, M., Stewart, J.S., Pedersen, S.H.,
 Reinking, A.K., Elder, K.: A lagrangian snow-evolution system for sea-ice applications (snowmodel-lg): Part i—model description. Journal of Geophysical Research:
 Oceans 125(10), 2019–015913 (2020)
- Liston, G.E., Mernild, S.H.: Greenland freshwater runoff. part i: A runoff routing
 model for glaciated and nonglaciated landscapes (hydroflow). Journal of Climate
 25(17), 5997–6014 (2012)
- Li, H., Sheffield, J., Wood, E.F.: Bias correction of monthly precipitation and temperature fields from intergovernmental panel on climate change ar4 models using
 equidistant quantile matching. Journal of Geophysical Research: Atmospheres
 115(D10) (2010)
- Naveau, P., Huser, R., Ribereau, P., Hannart, A.: Modeling jointly low, moderate, and
 heavy rainfall intensities without a threshold selection. Water Resources Research
 524 52(4), 2753–2769 (2016)

of the royal meteorological society 137(656), 553–597 (2011)

- Nash, J.E., Sutcliffe, J.V.: River flow forecasting through conceptual models part i—a
 discussion of principles. Journal of hydrology 10(3), 282–290 (1970)
- Olsson, J., Berggren, K., Olofsson, M., Viklander, M.: Applying climate model precipitation scenarios for urban hydrological assessment: A case study in kalmar city, sweden. Atmospheric Research **92**(3), 364–375 (2009)
- Papastathopoulos, I., Tawn, J.A.: Extended generalised pareto models for tail
 estimation. Journal of Statistical Planning and Inference 143(1), 131–143 (2013)
- Rasmussen, R., Liu, C., Ikeda, K., Gochis, D., Yates, D., Chen, F., Tewari, M., Barlage,
 M., Dudhia, J., Yu, W., *et al.*: High-resolution coupled climate runoff simulations
 of seasonal snowfall over colorado: A process study of current and warmer climate.
 Journal of Climate 24(12), 3015–3048 (2011)
- Schubert, S., Henderson-Sellers, A.: A statistical model to downscale local daily
 temperature extremes from synoptic-scale atmospheric circulation patterns in the
 australian region. Climate Dynamics 13, 223–234 (1997)
- Scarrott, C., MacDonald, A.: A review of extreme value threshold estimation and uncertainty quantification. REVSTAT-Statistical journal **10**(1), 33–60 (2012)
- Stein, M.L.: Parametric models for distributions when interest is in extremes with an
 application to daily temperature. Extremes 24, 293–323 (2021)
- Taylor, J.W.: A quantile regression neural network approach to estimating the conditional density of multiperiod returns. Journal of forecasting 19(4), 299–311 (2000)
- Teutschbein, C., Seibert, J.: Bias correction of regional climate model simulations
 for hydrological climate-change impact studies: Review and evaluation of different
 methods. Journal of hydrology 456, 12–29 (2012)
- ⁶⁰⁹ Vrac, M., Naveau, P.: Stochastic downscaling of precipitation: From dry events to
 ⁶¹⁰ heavy rainfalls. Water resources research 43(7) (2007)

Wang, L., Chen, W.: Equiratio cumulative distribution function matching as an
 improvement to the equidistant approach in bias correction of precipitation.
 Atmospheric Science Letters 15(1), 1–6 (2014)

Wilby, R.L., Hay, L.E., Leavesley, G.H.: A comparison of downscaled and raw gcm output: implications for climate change scenarios in the san juan river basin, colorado.
Journal of Hydrology 225(1), 67–91 (1999) https://doi.org/10.1016/S0022-1694(99)
00136-5

⁶¹⁸ Wood, A.W., Maurer, E.P., Kumar, A., Lettenmaier, D.P.: Long-range experimental ⁶¹⁹ hydrologic forecasting for the eastern united states. Journal of Geophysical Research:

 $_{620}$ Atmospheres **107**(D20), 6 (2002)