

# Impact of bias nonstationarity on the performance of uni- and multivariate bias-adjusting methods

Jorn Van de Velde<sup>1,2</sup>, Matthias Demuzere<sup>3,1</sup>, Bernard De Baets<sup>2</sup>, and Niko E. C. Verhoest<sup>1</sup>

<sup>1</sup>Hydro-Climatic Extremes Lab, Ghent University, Ghent, Belgium

<sup>2</sup>KERMIT, Department of Data Analysis and Mathematical Modelling, Ghent University, Ghent, Belgium

<sup>3</sup>Department of Geography, Ruhr-University Bochum, Bochum, Germany

**Correspondence:** Jorn Van de Velde (jorn.vandvelde@ugent.be)

## Abstract.

Climate change is one of the biggest challenges currently faced by society, with an impact on many systems, such as the hydrological cycle. To ~~locally~~ assess this impact in a local context, Regional Climate Model (RCM) simulations are often used as input for ~~hydrological~~ rainfall-runoff models. However, RCM results are still biased with respect to the observations.

5 Many methods have been developed to adjust these biases, but only during the last few years, methods to adjust biases that account for the correlation between the variables have been proposed. This correlation adjustment is especially important for compound event impact analysis. As ~~a simple example of those compound events, an illustration, a~~ hydrological impact assessment exercise is used here, as hydrological models often need multiple locally unbiased input variables to ensure an unbiased output. However, it has been suggested that multivariate bias-adjusting methods may perform poorly under climate

10 change conditions because of bias nonstationarity. In this study, two univariate and three multivariate bias-adjusting methods are compared with respect to their performance under climate change conditions. To this end, the methods are calibrated in the late 20th century (1970-1989) and validated in the early 21st century (1998-2017), in which the effect of climate change is already visible. The variables adjusted are precipitation, evaporation and temperature, of which the former two are used as input for a rainfall-runoff model, to allow for the validation of the methods on discharge. Although not used for discharge

15 modelling, temperature is a commonly-adjusted variable in both uni- and multivariate settings and ~~therefore important to take into account~~we therefore also included this variable in our research. The methods are also evaluated using indices based on the adjusted variables, the temporal structure, and the multivariate correlation. ~~For precipitation, all methods decrease the bias in a comparable manner~~The Perkins Skill Score is used to evaluate the full PDF. The results show a clear impact of nonstationarity on the bias adjustment. However, ~~for many other indices the results differ considerably between the bias-adjusting methods.~~

20 ~~The multivariate methods often perform worse than the univariate methods, a result that is especially notable for temperature and evaporation. As these variables have already changed the most under climate change conditions, this reinforces the opinion that the~~the impact varies depending on season and variable: the impact is most visible for precipitation in winter and summer. This should be accounted for in both multivariate bias-adjusting methods ~~are not yet fit to cope with nonstationary climate conditions. Although the effect is slightly dampened by the hydrological model, our analysis still reveals that, to date, the~~

25 ~~simpler univariate bias-adjusting methods are preferred for assessing climate change impact and impact models. In the former~~

because these do not always include seasonality; in the latter because incorrectly adjusted inputs or forcings will lead to predicted discharges that are biased.

Copyright statement. TEXT

## 1 Introduction

30 The influence of climate change is felt throughout many regions of the world, as becomes evident from the higher frequency or intensity of natural hazards, such as floods, droughts, heatwaves and forest fires (IPCC, 2012). As these intensified natural hazards threaten society, it is essential to be prepared for them. Knowledge on future climate change is obtained by running Global Climate Models (GCMs), creating large ensemble outputs such as in the Climate Model Intercomparison Project 6 (CMIP6) (Eyring et al., 2016). Although they are informative on a global scale, the generated data are too coarse for local  
35 climate change impact assessments. To bridge the gap from the global to the local scale, Regional Climate Models have become a standard application (Jacob et al., 2014), using the output from GCMs as input or boundary conditions.

Although the information provided by both GCMs and RCMs is ~~very~~-valuable, both are biased ~~with respect to~~ w.r.t. the observations, especially for precipitation (Kotlarski et al., 2014). The biases can occur in any statistic and are commonly defined as *“a systematic difference between a simulated climate statistic and the corresponding real-world climate statistic”* (Maraun, 2016). These biases are caused by temporal or spatial discretisation and unresolved or unrepresented physical processes (Teutschbein and Seibert, 2012; Cannon, 2016). An important example of the latter is convective precipitation, which can only  
40 be resolved by very high resolution models. Although the further improvement of models is an important area of research (Prein et al., 2015; Kendon et al., 2017; Helsen et al., 2019; Fosser et al., 2020), such improved models are computationally expensive. As such, it is still necessary practice to statistically adapt the climate model output to adjust the biases (Christensen  
45 et al., 2008; Teutschbein and Seibert, 2012; Maraun, 2016).

Many different bias-adjusting methods exist (Teutschbein and Seibert, 2012; Gutiérrez et al., 2019). They all calibrate a transfer function using the historical simulations and historical observations and apply this transfer function to the future simulations to generate future ‘observed values’ or an adjusted future. Of all the different methods, the quantile mapping method (Panofsky et al., 1958) was shown to be the generally best performing method (Rojas et al., 2011; Gudmundsson et al.,  
50 2012). Quantile mapping adjusts biases in the full distribution, whereas most other methods only adjust biases in the mean and/or variance.

An important problem with quantile mapping and most other commonly used methods is that they are univariate and do not adjust biases in the multivariate correlation. Although quantile mapping can retain climate model multivariate correlation (Wilcke et al., 2013), the ability of univariate methods to improve the climate model’s multivariate correlation has been ques-  
55 tioned (Hagemann et al., 2011; Ehret et al., 2012; Hewitson et al., 2014). This is important for impact assessment, as local impact models often need multiple input variables and many high-impact events are caused by the co-occurrence of multiple

phenomena, the so-called ‘compound events’ (Zscheischler et al., 2018, 2020). For example, ~~floods can be characterised~~ flood magnitude can be projected by a rainfall-runoff model using evaporation and precipitation time series as an input. If the correlation between these variables is biased ~~with respect to~~ w.r.t. the observations, then it can be expected that the model output is biased as well. ~~This results in a higher uncertainty when using these models and thus in the resulting assessment,~~ which can further propagate in the impact models. During the past decade, multiple methods have been developed to counter this problem. The first methods focused on the adjustment of two jointly occurring variables, most often precipitation and temperature, such as those by Piani and Haerter (2012) and Li et al. (2014). However, it became clear that adjusting only two variables would not suffice, hence many more methods have been developed that jointly adjust ~~more~~ multiple variables, including those by Vrac and Friederichs (2015); Cannon (2016); Mehrotra and Sharma (2016); Dekens et al. (2017); Cannon (2018); Vrac (2018); Nguyen et al. (2018); Robin et al. (2019). Yet, the recent growth in availability of such methods comes along with a gap in the knowledge on their performance. In some studies, these methods have been compared with one or two older multivariate methods to reveal the improvements (Vrac and Friederichs, 2015; Cannon, 2018) or with univariate methods (Räty et al., 2018; Zscheischler et al., 2019; Meyer et al., 2019). Each of ~~these three studies indicates that the univariate and multivariate methods~~ the latter three studies comparing uni- and multivariate bias adjusting methods indicates that these lead to different results, yet it is difficult to conclude whether uni- or multivariate methods perform best. According to Zscheischler et al. (2019) multivariate methods have an added value. Räty et al. (2018) conclude that the multivariate methods and univariate methods ~~performed~~ perform similarly, while Meyer et al. (2019) could not draw definitive conclusions. These studies vary in set-up, adjusted variables and study area, which all could have caused the difference in added value. In all three studies, the same method, namely the Multivariate Bias Correction in  $n$  dimensions (MBCn) (Cannon, 2018) was the basis for comparison. Only recently, the first studies comparing multiple multivariate bias-adjusting methods were published (François et al., 2020; Guo et al., 2020). The study by François et al. (2020) focused on the different principles underlying the multivariate bias-adjusting methods and concluded that the choice of method should be based on the end user’s goal. Besides, they also noticed that ~~so far,~~ all multivariate methods ~~fail in representing~~ studied fail in adjusting the temporal structure of a time series. In contrast to the focus of François et al. (2020), Guo et al. (2020) studied the performance of multivariate bias-adjusting methods for climate change impact assessment and concluded that multivariate methods could be interesting in this context. However, they also noticed that the performance of the multivariate methods was lower in the ~~more recent~~ validation period and suggested that this could be caused by bias nonstationarity. As the use of multivariate bias-adjusting methods could be an important tool for climate change impact assessment, this deserves more attention.

The bias stationarity - or bias time invariance - assumption is the most important assumption for bias correction. It implies that the bias is the same in the calibration and validation or future periods and that the transfer function based on the calibration period can ~~consequently~~ thus be used in the future period. However, this assumption does not hold due to different types of nonstationarity induced by climate change, which may cause problems (Milly et al., 2008; Derbyshire, 2017). In the context of bias adjustment, this problem has been known for several years (Christensen et al., 2008; Ehret et al., 2012), but has not received a lot of attention. A few authors have tried to propose new types of bias relationships (Buser et al., 2009; Ho et al., 2012; Sunyer et al., 2014; Kerkhoff et al., 2014). Recently, it has been suggested that it is best to assume a non-monotonic

bias change (Van Schaeybroeck and Vannitsem, 2016). Some authors suggested that bias nonstationarity could be an important source of uncertainty (Chen et al., 2015; Velázquez et al., 2015; Wang et al., 2018; Hui et al., 2019), but not all found clear indications of bias nonstationarity (Maraun, 2012; Piani et al., 2010; Maurer et al., 2013).

95 The availability of new methods and more data enables a more coherent assessment of the bias (non)stationarity issue. By comparing ~~three-four~~ bias-adjusting methods in a climate change context with possible bias nonstationarity, some of the remaining questions in François et al. (2020) and Guo et al. (2020) can be answered. The ~~three-four~~ multivariate bias-adjusting methods ~~that will be~~ compared in this study are ‘Multivariate Recursive Quantile Nesting Bias Correction’ (MRQNBC, Mehrotra and Sharma (2016)), MBCn (Cannon, 2018) ~~and~~, ‘dynamical Optimal Transport Correction’ (dOTC, Robin et al. (2019))  
100 ~~-These three and~~ ‘Rank Resampling for Distributions and Dependences’ (R<sup>2</sup>D<sup>2</sup>, Vrac (2018); Vrac and Thao (2020b)). ~~These four~~ methods give a broad view of the different multivariate bias adjustment principles, which we will elaborate on in Section 3.3. As a baseline, two univariate bias-adjusting methods will be used: Quantile Delta Mapping (QDM, Cannon et al. (2015)) and modified Quantile Delta Mapping (mQDM, Pham (2016)). QDM is a classical univariate bias-adjusting method and is chosen for this analysis as it is a robust and relatively common quantile mapping method, especially as one of the sub-  
105 routines in the multivariate bias-adjusting methods (Mehrotra and Sharma, 2016; Nguyen et al., 2016; Cannon, 2018). mQDM, on the other hand, is one of the so-called ‘delta change’ methods, which are based on an adjustment of the historical time series. Using these univariate bias-adjusting methods, we can assess whether multivariate and univariate bias-adjusting methods differ in their response to possible bias nonstationarity.

The methods will be compared by applying them for the bias adjustment of precipitation, potential evaporation and tempera-  
110 ture. The bias-adjusted time series will be used as inputs for a hydrological model in order to simulate the discharge. Discharge time series are the basis for flood hazard calculation, but can also be considered as an interesting source of validation themselves (Hakala et al., 2018). ~~Although temperature is not needed as an input for~~ The bias adjustment and discharge simulation are both assessed at one grid cell/location only. Although this does not allow for investigating the spatial extent and impact of nonstationarity, the hydrological model, it is, together with precipitation, the most common variable to be adjusted in similar  
115 ~~studies and therefore it is also included here. In order to mimic climate change context, the ‘historical’ or calibration time series runs from 1970 to 1989 and the ‘future’ or validation time series runs from 1998 to 2017, which is only recent past. In the latter time frame, effects of climate change are already visible (IPCC, 2013)~~ focus on one location gives information on the influence of possible bias nonstationarity on local impact models and may hence be a starting point for broader assessments. We will also not account for the differences between models, as we only investigate a single GCM-RCM model chain. This  
120 allows for a precise investigation of the possible effects of bias nonstationarity, although it does not allow for assessing other types of uncertainty. The change of some biases from calibration to validation time series will be calculated, to indicate the extent of the bias nonstationarity. Maurer et al. (2013) proposed the R index for this purpose (~~see Section 2.4~~). Calculating the bias nonstationarity between both periods will give an indication of the impact of a changing bias on climate impact studies for the end of the 21st century. As Chen et al. (2015) mentioned: *“If biases are not constant over two very close time periods,*  
125 *there is little hope they will be stationary for periods separated by 50 to 100 years”*

## 2 Data and validation

### 2.1 Data

~~For the observations, the dataset made available by the~~ The observational data used were obtained from the Belgian Royal Meteorological Institute ~~and described in Van de Velde et al. (2020) is also used in this study. This dataset comprises (RMI)~~  
130 ~~Uccle observatory. The most important time series used is the 10-min precipitation amount, gauged with a Hellmann-Fuess pluviograph, from 1898 to 2018. An earlier version of this precipitation dataset was described by Demarée (2003) and analyzed in De Jongh et al. (2006). Multiple other studies have used this time series (Verhoest et al., 1997; Verstraeten et al., 2006; Vandenberghe et al., 2006). The 10-min precipitation time series was aggregated to daily level to be comparable with the other time series used.~~

~~For the multivariate methods, the precipitation time series was combined with a 2 meter air temperature and potential~~  
135 ~~evaporation time series. The daily potential evaporation was calculated by the RMI from 1901 to 2019, using the Penman formula for a grass reference surface (Penman, 1948) with variables measured at the Uccle observatory. Daily average temperatures were obtained using measurements from 1901 to 2019. As the last complete year for precipitation was 2017, the data were used from 1901 to 2017, amounting to 117 years (1901-2017) of daily precipitation amount, daily average temperature and daily potential evaporation years of daily data. As Uccle (near Brussels) is situated in a region with small topographic differences, it~~  
140 ~~is assumed that the precipitation statistics within the grid cell are uniform. Hence, the Uccle data can be used for comparison with the gridded climate simulation data discussed below.~~

~~The IPCC report (IPCC, 2013) clearly states the influence of climate change on different variables. For Belgium, this is illustrated by Fig. ??, in which the temperature and evaporation anomalies for the 21st century are all higher than the long-term mean value. However, for precipitation, the effect of climate change is not yet visible.~~

145 ~~Yearly mean temperature, precipitation and evaporation anomalies for 1901-2017, compared with long-term mean value from 1920-1980. Red points are 21st century values.~~

~~As in Van de Velde et al. (2020),~~ For the simulations, data from the EURO-CORDEX (Jacob et al., 2014) project (Jacob et al., 2014) were used. The Rossby Centre regional climate model RCA4 was used (Strandberg et al., 2015), with MPI-ESM-LR GCM (Popke et al., 2013) boundary conditions. RCA4 is used (Strandberg et al., 2015) as it is one of the few RCMs with potential evaporation as an output variable. This RCM was forced with boundary conditions from the MPI-ESM-LR GCM (Popke et al., 2013) and has a spatial resolution of 0.11°, or 12.5 km. Historical data and scenario data for the grid cell comprising Uccle were respectively obtained for 1970-2005 and 2006-2100. The former time frame is limited by the earliest available data from the RCM. The latter time frame was only used until 2017, in accordance with the observational data. As climate change scenario, an RCP4.5 forcing was used in this paper (van Vuuren et al., 2011). Since only 'near future' (from  
155 ~~the model point of view) data were used, the choice of forcing does not have a large impact. However, when studying scenarios in a time frame further away from the present, using an ensemble of forcings is more relevant to be aware of the uncertainty regarding future climate change impact. evaluations of the RCA4 model have shown that there is a bias in precipitation, especially in winter (Strandberg et al., 2015), but this bias is in line with the biases from other EURO-CORDEX models (Kotlarski et al., 2014).~~

## 160 2.2 Time frames

As mentioned in the introduction, it is important to assess bias-adjusting methods in a context they will be used in, i.e. under climate change conditions. The time series used in this study were chosen accordingly: 1970-1989 was chosen as the ‘historical’ or calibration time period and 1998-2017 was chosen as the ‘future’ or validation time period. [In this time frame, effects of climate change are already visible \(IPCC, 2013\)](#). Time series of 20 years were chosen here, although it is advised to use 30 years of data to have robust calculations (Berg et al., 2012; Reiter et al., 2018). However, as no climate model data prior to 1970 are available, using 30 years of data would have led to overlapping time series.

## 2.3 Validation framework

An important aspect in bias adjustment is the validation of the methods. Different methods are available, of which a pseudo-reality experiment (Maraun, 2012) is one of the most-used ones. In this method, each member of a model ensemble is in turn used as the reference in a cross-validation. However, while such a set-up is useful when comparing bias-adjustment methods, it only mimics a real application context. When sufficient observations are available, a ‘pseudo-projection’ [setup set-up](#) (Li et al., 2010) can be used. This set-up resembles a ‘differential split-sample testing’ (Klemeš, 1986) and is more in agreement with a practical application of bias-adjusting methods. Differential split-sample testing has been used in a bias adjustment context by Teutschbein and Seibert (2013), by constructing two time series with respectively the driest and wettest years. In our case study, it is assumed that the two time series differ enough because of climate change. Consequently, the approach is simple, and as the validation is not set in the future, it is considered a ‘pseudo-projection’.

Besides the choice of time frames and data, also the choice of validation indices is of key importance. Maraun and Widmann (2018a) stress that these indices should only be indirectly affected by the bias adjustment, as only validating on adjusted indices can be misleading. Such adjusted indices are the precipitation intensity, temperature and evaporation, which are used to build the transfer function in the historical setting and should be corrected by construction. Under bias stationarity, this correction will be carried over to the future, possibly hiding small inconsistencies that may arise for extreme values. If the bias is not stationary, the effect might be different between adjusted and indirectly affected indices. As such, besides the three adjusted variables (indices 1 to 3 in Table 1) and their correlations (indices 4 to 12, which are directly adjusted by some of the methods), also indices based on the precipitation occurrence and on the discharge  $Q$  are used. The occurrence-based indices (13 to 16) allow for assessing how the methods influence the precipitation time series structure, ~~an influence that might be potentially large (Van de Velde et al., 2020)~~. The discharge-based indices (17 and 18) allow for the assessment of the impact of the different bias-adjusting methods on simulated river flow. The discharge-based indices combine the information of the other indices by routing through the rainfall-runoff model. They are the most important aspect of the assessment, as they indicate the natural hazard. ~~ETCCDI (Expert Team on Climate Change Detection and Indices) precipitation indices (Zhang et al., 2011) have also been considered and calculated. However, these are not included in this paper, as the differences in ETCCDI indices were minor and did not allow to clearly discern between the different methods~~ [As the percentiles focus mostly on the extremes, the Perkins Skill Score \(PSS\)](#) (Perkins et al., 2007) [is used to assess the adjustment of the full PDF of the variables](#). All indices

were~~are~~ calculated taking all days into account, instead of only calculating them on wet days, as some of the multivariate bias-adjusting methods do not discriminate between wet or dry days in their adjustment.

195 The indices are all calculated on a seasonal basis for both the calibration and validation period. By comparing over these periods, we can relate the performance to either the method itself or bias (non)stationarity, on a seasonal basis. Besides, not all methods adjust on a seasonal basis. As such, methods performing poorly in both periods might need a seasonal component for bias adjustment. The seasons were defined as follows: winter (DJF), spring (MAM), summer (JJA) and autumn (SON).

**Table 1.** Overview of the indices used

| Nr | Index                      | Name   |
|----|----------------------------|--|
| 1  | $P_x$                      | Precipitation amount percentile values, with $x$ the percentile considered |
| 2  | $T_x$                      | Temperature percentile values, with $x$ the percentile considered          |
| 3  | $E_x$                      | Evaporation percentile values, with $x$ the percentile considered          |
| 4  | $\text{corr}_{P,E}$        | Spearman correlation between the time series of $P$ and $E$                |
| 5  | $\text{corr}_{P,T}$        | Spearman correlation between the time series of $P$ and $T$                |
| 6  | $\text{corr}_{E,T}$        | Spearman correlation between the time series of $E$ and $T$                |
| 7  | $\text{crosscorr}_{P,E,0}$ | Lag-0 crosscorrelation between the time series of $P$ and $E$              |
| 8  | $\text{crosscorr}_{P,T,0}$ | Lag-0 crosscorrelation between the time series of $P$ and $T$              |
| 9  | $\text{crosscorr}_{E,T,0}$ | Lag-0 crosscorrelation between the time series of $E$ and $T$              |
| 10 | $\text{crosscorr}_{P,E,1}$ | Lag-1 crosscorrelation between the time series of $P$ and $E$              |
| 11 | $\text{crosscorr}_{P,T,1}$ | Lag-1 crosscorrelation between the time series of $P$ and $T$              |
| 12 | $\text{crosscorr}_{E,T,1}$ | Lag-1 crosscorrelation between the time series of $E$ and $T$              |
| 13 | $P_{P00}$                  | Precipitation transition probability from a dry to a dry day               |
| 14 | $P_{P10}$                  | Precipitation transition probability from a wet to a dry day               |
| 15 | $N_{\text{dry}}$           | Number of dry days   |
| 16 | $P_{\text{lag}1}$          | Precipitation lag-1 auto-correlation                                       |
| 17 | $Q_x$                      | Discharge percentiles, with $x$ the percentile considered                  |
| 18 | $Q_{T20}$                  | 20-year return period value of discharge                                   |

## 2.4 Bias nonstationarity

200 In a study on possible changes in bias, Maurer et al. (2013) proposed the R index:

$$R = 2 \frac{|\text{bias}_f - \text{bias}_h|}{|\text{bias}_f| + |\text{bias}_h|}, \quad (1)$$

where  $\text{bias}_f$  and  $\text{bias}_h$  are the biases in respectively the future and historical time series, calculated on the basis of the observations and raw climate simulations. The R index takes a value between 0 and 2. If the index is greater than one, the difference in bias between the two sets is larger than the average bias of the model and it is likely that the bias adjustment would degrade

205 the RCM output rather than improve it. The index is calculated for the indices used for validation in order to have an indication of the influence of bias nonstationarity on these indices. Besides for the indices, the R index is also calculated for the average and standard deviation of each variable, in order to be able to more easily visualise the changes in distribution.

## 2.5 Hydrological model

Similar to ~~Van de Velde et al. (2020)~~Pham et al. (2018), we use the Probability Distributed Model (PDM, ~~Moore (2007)~~Moore (2007); Cabus  
210 a lumped conceptual rainfall-runoff model to calculate the discharge for the Grote Nete watershed in Belgium. This model uses precipitation and evaporation time series as inputs to generate a discharge time series. The PDM as used here was calibrated by ~~Cabus (2008)~~ (RMSE = 0.9 m<sup>3</sup>/h, see Pham et al. (2018) for more details) using the Particle Swarm Optimization algorithm (PSO, ~~Eberhart and Kennedy (1995)~~). ~~The same assumption was used as in Pham et al. (2018) and Van de Velde et al. (2020), i.e.~~Eberhart and Kennedy (1995)). As in Pham et al. (2018), it was assumed that the differences between meteorological conditions in the Grote Nete-watershed and Uccle are negligible, and ~~thus that that thus~~ the adjusted data for the Uccle grid cell  
215 can be used as a forcing for the PDM. This assumption is based on the limited distance of 50 km between the gauging stations used for the observations in Uccle and the gauging station used for the PDM calibration. As mentioned before, the region has a flat topography and, hence, the climatology can be considered similar. Furthermore, the goal is not to make predictions, but to assess the impact of different bias adjustment methods on the discharge values. To calculate the bias on the ~~discharge~~  
220 ~~indices, both the indices, observed,~~ raw and adjusted ~~precipitation and evaporation~~ RCM time series were used as forcing for this model. The discharge time series generated by the observations is considered to be the 'observed' discharge, and biases are calculated in comparison with this time series.

## 2.6 Validation metrics

The residual biases relative to the observations and to the model bias are often used in this paper to graphically present and  
225 interpret the results(~~Van de Velde et al., 2020~~). These residual biases are based on the 'added value' concept (Di Luca et al., 2015) and enable a comparison based on two aspects. The first aspect is the performance in removing the bias, the second is the extent of the bias removal in comparison with the original value for the corresponding index for the observation time series. The use of the residual biases allows for a detailed study and comparison of the effect of bias adjustment on the different indices.

230 The residual bias relative to the observations  $RB_O$  for an index  $k$  is calculated as follows:

$$RB_O(k) = 1 - \frac{|\text{bias}_{\text{raw}(k)}| - |\text{bias}_{\text{adj}(k)}|}{|\text{obs}(k)|}, \quad (2)$$

with  $\text{raw}(k)$  the raw climate model simulations,  $\text{adj}(k)$  the adjusted climate model simulations and  $\text{obs}(k)$  the observed values for index  $k$ .

The residual bias relative to the model bias  $RB_{MB}$  for an index  $k$  is calculated as follows:

235  $RB_{MB}(k) = 1 - \frac{|\text{bias}_{\text{raw}(k)}| - |\text{bias}_{\text{adj}(k)}|}{|\text{bias}_{\text{raw}(k)}|}.$  (3)



Absolute values are used in Eqs. (2) and (3) to compute the absolute difference between the raw and adjusted values, thus neglecting a possible change of sign of the bias. If the values of these residual biases are lower than 1 for an index, the method performs better than the raw RCM for this index. The best methods have low scores on both residual biases for as many indices as possible.

## 240 3 Bias-adjusting methods

### 3.1 Occurrence-bias adjustment: Thresholding

One of the deficiencies of RCMs, especially in Northwest Europe, are the so-called ‘drizzle days’ (Gutowski et al., 2003; Themeßl et al., 2012; Argüeso et al., 2013), ~~i.e. the simulation of a small amount of precipitation on days that are supposed to be during which small amounts of precipitation are simulated while these days should have been~~ dry. This has an influence on the temporal structure of the simulated time series and should thus be adjusted (Ines and Hansen, 2006). This is commonly done in an occurrence-bias-adjusting step before the main step, the intensity-bias adjustment. ~~Although multiple methods are available (Hay and Clark, 2003; Schmidli et al., 2006; Vrae et al., 2016; Van de Velde et al., 2020), their effect on the intensity-bias-adjusting methods is not always clear, especially if these add a lot of complexity, such as the multivariate intensity-bias-adjusting methods. Furthermore, the more advanced methods, with stochastic steps, do not always seem to have~~ an added value (Van de Velde et al., 2020). Therefore, we proposed to use simpler methods, such as thresholding .

~~Thresholding~~ In this study, we use the thresholding occurrence-bias-adjusting method, which is one of the most common occurrence-bias-adjusting methods ~~and has been in use for many years~~ (e.g. Hay and Clark (2003); Schmidli et al. (2006); Ines and Hansen (2006)). This method is only applicable in regions where the assumption holds that the simulated time series has more wet days than the observed time series. This is the case for Northwest Europe (Themeßl et al., 2012) and Belgium in particular. An advanced version of the thresholding method is used here. To adjust the number of wet days, the ~~frequencies~~ total number of dry days in the observations and in the simulations are calculated. The difference in dry days between the two ~~frequencies~~ periods,  $\Delta N$ , is the number of days of the simulated time series that have to be adapted. ~~The simulated series is adapted by first sorting the wet days and thus only changing~~ If  $\Delta N$  days have to be converted to dry days, then the  $\Delta N$  lowest days of the simulation time series by setting them to 0. ~~days with the lowest amounts of precipitation are changed to dry days.~~

$\Delta N$  is computed for the past and applied in the future and consequently relies on the bias stationarity assumption. However, as thresholding is used prior to all methods, the influence of possible bias nonstationarity on  $\Delta N$  is assumed to be negligible. ~~More mathematical details and an algorithm including all steps can be found in Van de Velde et al. (2020). In the present paper, days in the simulation are~~ Besides, as is shown in Section 4.1, the number of dry days is stationary for the time frames studied in this paper.

In this advanced version of thresholding, some considerations are made. First, a day is considered wet if ~~the daily precipitation intensity is higher than its simulated precipitation amount is above~~ 0.1 mm, to account for measurement errors in the observations. Second, the adjustment is done on a monthly basis, to account for the temporal structure in the observed time series.

Third, both historical and future simulations are adjusted, to ensure that the bias can be transferred from the historical to the future time period would be impaired.

## 270 3.2 Univariate intensity-bias-adjusting methods

### 3.2.1 Quantile Delta Mapping

The Quantile Delta Mapping (QDM) method was first proposed by Li et al. (2010). Its main idea is to preserve the climate simulation trends: it takes trend nonstationarity (changes in the simulated distribution) into account to a certain degree. Although it handles temperature adjustments well, it gives unrealistic values for precipitation and was therefore extended by Wang and Chen (2014) for precipitation adjustment. ~~A comparison with other quantile mapping methods by Cannon et al. (2015) showed this method to perform best with respect to the preservation of trends. Cannon et al. (2015) bundled both the method by~~ By combining the methods by Li et al. (2010) (*Equidistant CDF-matching*) and Wang and Chen (2014) (*Equiratio CDF-matching*) ~~under the name~~ Quantile Delta Mapping, because of the similarity with delta change methods (which are described in e. g. Olsson et al. (2009), Willem's and Vrae (2011) and Rätty et al. (2014))., Cannon et al. (2015) developed the Quantile Delta Mapping method.

Mathematically, this method can be written as

$$x_i^{\text{fa}} = x_i^{\text{fs}} + F_{x^{\text{ho}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}})) - F_{x^{\text{hs}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}})) \quad (4)$$

in the additive case, and

$$x_i^{\text{fa}} = x_i^{\text{fs}} \frac{F_{x^{\text{ho}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}}))}{F_{x^{\text{hs}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}}))} \quad (5)$$

285 in the ratio or multiplicative case. The superscripts hs, ho, fs and fa indicate respectively the historical simulations, the historical observations, the future simulations and the adjusted future. ~~The~~ In this paper, the additive version is used for temperature adjustments, ~~whereas the multiplicative version is used for the adjustment of precipitation and evaporation, to ensure physically correct values (Hempel et al., 2013). The same implementation is used as in Van de Velde et al. (2020): a 91-day window to time series and the multiplicative one for precipitation and evaporation time series. This choice is based on the~~ work of Wang and Chen (2014), who have shown that using the additive adjustment for precipitation results in unrealistic precipitation values and introduced a multiplicative adjustment. For evaporation, we follow the few available studies (e.g. Lenderink et al. (2007)) in using the same adjustment as for precipitation.

To ensure the consistency of the time series (Thiemeßl et al., 2011; Rajczak et al., 2016; Reiter et al., 2018); empirical CDFs for the ease of application (Gudmundsson et al., 2012; Gutjahr and Heinemann, 2013); application of the multiplicative version on wet days only, a 91-day moving window is opted for, as suggested by Rajczak et al. (2016) and Reiter et al. (2018). This enables the adjustment of each day based on 91 days/year · 20 years = 1820 days. These days were used to build an empirical CDF (as in Gudmundsson et al. (2012); Gutjahr and Heinemann (2013), among others). It is also important to note that for precipitation, Eq. (5) was applied only on the days considered wet, i.e. with a precipitation /evaporation higher than 0.1 mm.

For consistency, a threshold of 0.1 mm /day- was also used for evaporation. It is important to note that although QDM is  
 300 only applied on wet days, it can still transform low-precipitation wet days into days that are considered to be dry (e.g. with a  
 precipitation amount < 0.1 mm) if the ratio in Eq. (5) is small enough.

### 3.2.2 Modified Quantile Delta Mapping

Pham (2016) proposed another version of QDM, following the delta change philosophy (Olsson et al., 2009; Willems and Vrac,  
 2011): the trend established by the RCM is assumed to be more thrust-worthy than the absolute value itself. When applying  
 305 this type of methods, the simulated change between the historical and the future is applied to the observations. Thus, instead  
 of the future simulations, the historical observations are adjusted to the future ‘observations’. As Johnson and Sharma (2011)  
 mention, this workflow could be problematic for future impact assessment, as it inherits the temporal structure of the historical  
 observations. This method is mathematically very similar to the QDM method, exchanging the roles of  $x^{\text{fs}}$  and  $x^{\text{ho}}$ . Thus, it is  
 named ‘modified Quantile Delta Mapping’ (mQDM), and can for the additive case be written as

$$310 \quad x_i^{\text{fa}} = x_i^{\text{ho}} + F_{x^{\text{fs}}}^{-1} \left( F_{x^{\text{ho}}} \left( x^{\text{ho}} \right) \right) - F_{x^{\text{hs}}}^{-1} \left( F_{x^{\text{ho}}} \left( x^{\text{ho}} \right) \right). \quad (6)$$

The ratio version is ~~mathematically written as~~ given by

$$x_i^{\text{fa}} = x_i^{\text{ho}} \frac{F_{x^{\text{fs}}}^{-1} \left( F_{x^{\text{ho}}} \left( x^{\text{ho}} \right) \right)}{F_{x^{\text{hs}}}^{-1} \left( F_{x^{\text{ho}}} \left( x^{\text{ho}} \right) \right)}. \quad (7)$$

For the implementation, the same principles were used as for the QDM method: a 91-day moving window, empirical CDFs  
 and a ~~threshold~~ minimum value of 0.1 mm/day to be considered as a wet day.

### 315 3.3 Multivariate intensity-bias-adjusting methods

The increasing number of multivariate bias-adjusting methods throughout the 2010s urges the need to classify them according  
 to their properties. One possible classification was done by Vrac (2018), who proposed the ‘marginal/dependence’ versus the  
 ‘successive conditional’ approach. The former approach separately adjusts the 1D-marginal distributions and the dependence  
 structure and is applied in e.g. Vrac and Friederichs (2015), Cannon (2018) and Vrac (2018). These two components are then  
 320 recombined to obtain data that are close to the observations for both marginal and multivariate aspects. The latter approach  
 consists of adjusting ~~one given variable and then adjusting a second a~~ variable conditionally on the ~~second variable: this~~  
~~variables already adjusted. This~~ procedure is applied successively to each variable. Examples can be found in e.g. Piani and  
 Haerter (2012), Li et al. (2014) and Dekens et al. (2017). ~~According to Vrac (2018), the latter approach suffers from two~~  
~~main limitations. First, the adjustment is performed conditionally on the previously adjusted data. However, the adjustment~~  
 325 ~~is often applied in bins. As a result, for each variable, the amount of data available for each bin decreases, thus decreasing~~  
~~the robustness of the adjustment. Second, the ordering of the variables in the successive adjustments matters. For example,~~  
 Li et al. (2014) point out that their ‘Joint Bias Correction for temperature’ (JBCt) and ‘Joint Bias Correction for precipitation’  
 (JBCp) methods, which respectively first adjust temperature and precipitation, differ in performance. For these two reasons,

~~Vrac (2018)~~ Vrac (2018) discusses the limitations of the ‘successive conditional’ approach and advocates for the use of the more robust and coherent ‘marginal/dependence’ approach. Hence, ‘successive conditional’ methods are not included in the present paper. Robin et al. (2019) and François et al. (2020) extended ~~this the~~ classification by introducing the ‘all-in-one’ approach, which adjusts the marginal variables and the correlations simultaneously, ‘dynamical Optimal Transport Correction’ (dOTC) (Robin et al., 2019) being such a method.

Another perspective on the multivariate bias-adjusting methods is to consider the amount of temporal adjustment that is allowed or applied by the method. This is important, as the amount of temporal adjustment is intrinsically linked with the main goal, the adjustment of the multivariate distribution of the variables. This distribution, in which the dependence is characterised by the underlying copula (Nelsen, 2006; Schölzel and Friederichs, 2008), can be estimated using the ranks. Thus, to adjust the multivariate distribution, the ranks of the climate model are replaced by those of the observations, using methods such as the ‘Schaake Shuffle’ (Clark et al., 2004; Vrac and Friederichs, 2015). This implies that the temporal structure and trends of the climate model will be altered, which may have a considerable impact (Van de Velde et al., 2020; François et al., 2020) (François et al., 2020). This impact is especially large when multiday characteristics strongly matter, such as in applications as the hydrological example we use in this study (Addor and Seibert, 2014). Vrac (2018) mentions this necessity to modify the temporal structure and rank chronology of the simulations. Yet, he also mentions that the extent of this modification is still a matter of debate. Cannon (2016) describes this as the ‘knobs’ that control whether marginal distributions, inter-variable or spatial dependence structure and temporal structure are more informed by the climate model or the observations. Thus, the choice between the temporal structure of the climate model and unbiased dependence structures is a trade-off that has to be made. Some methods, such as those by Vrac and Friederichs (2015), Mehrotra and Sharma (2016) and Nguyen et al. (2018) rely on the observations for their temporal properties, while other methods try to find the middle ground (e.g. Vrac (2018) and Cannon (2018)). A last perspective, which is not limited to multivariate methods, is that of trend preservation, i.e., the capacity of methods to preserve the changes simulated by the climate model, such as changes in mean, extremes and temporal structure. Although the amount of trend preservation or adjustment has been a matter of debate (Ivanov et al., 2018), Maraun (2016) argues that it is sensible to preserve the simulated changes and hence the climate change signal, if the model simulation is credible. As such, trend preservation interacts with bias nonstationarity: non-stationarity can be seen as the divergence between the observed and simulated trends. Hence, in a nonstationary context, trend-preserving methods may be disadvantaged, as they will assume that the simulated trend is trustworthy. In the univariate setting, QDM is an example of a trend-preserving method, as is ‘Scaled Distribution Mapping’ by Switanek et al. (2017).

Our choice of multivariate bias-adjusting methods takes the above classification into account. The oldest method in the comparison is ‘Multivariate Recursive Quantile Nesting Bias Correction’ (MRQNBC) (Mehrotra and Sharma, 2016). This method ~~completely~~ replaces the simulated correlations by those of the observations and is a ‘marginal/dependence’ method according to François et al. (2020). As QDM is used for the marginal distributions, the latter are preserved. However, MRQNBC does not preserve the changes in dependence. ‘Multivariate Bias Correction in  $n$  dimensions’ (Cannon, 2018) is both a ‘marginal/dependence’ method and a method that tries to combine information from the climate model and the observations. ~~The~~ Similar to MRQNBC, it explicitly preserves the simulated changes in the marginal distributions by applying QDM for the

marginal distributions. As the simulated dependence structure is the basis for the adjustment, it will be slightly preserved. The  
365 ‘Rank Resampling for Distributions and Dependences’ ( $R^2D^2$ , Vrac (2018); Vrac and Thao (2020b)) method preserves the  
rank correlation of the observations, but allows the climate model to have some influence on the temporal properties. It is also  
a ‘marginal/dependence’ method: in the present paper, QDM is used as its univariate routine and thus the changes in marginal  
distributions are preserved by  $R^2D^2$ . The most recent method, ‘dynamical Optimal Transport Correction’ (Robin et al., 2019)  
differs considerably from the other two methods: it generalises the ‘transfer function’-principle using the ‘optimal transport’  
370 paradigm (Villani, 2008), thereby defining a new category of multivariate bias-adjusting methods: the above-mentioned all-  
in-one approach. It is the only method that explicitly preserves the simulated changes in both the marginal distributions and  
the dependence structure. Although far from complete, by comparing these ~~three~~ four methods, a broad view of the different  
approaches in multivariate bias adjustment can be obtained. The main principles of the bias-adjusting methods are summarized  
in Table 2.

### 375 3.3.1 Multivariate Recursive Quantile Nesting Bias Correction

In 2016, Mehrotra and Sharma proposed a new multivariate bias adjustment method, named ‘Multivariate Recursive Quantile  
Nesting Bias Correction’ (MRQNBC), based on a combination of several older methods by Johnson and Sharma (2012),  
Mehrotra and Sharma (2012) and Mehrotra and Sharma (2015) and by incorporating QDM as the univariate routine for  
adjusting the marginals. The underlying idea of this method is to adjust on more than one timescale, ~~an idea that most  
380 bias-adjusting methods do not incorporate (Haerter et al., 2011). This~~ and to nest the results of the different timescales within  
each other. The adjustment on multiple timescales is ~~applied by adjusting almost never incorporated in bias-adjusting methods  
(Haerter et al., 2011). On each timescale,~~ the biases in lag-0- and lag-1-auto and the cross-correlation coefficients, i.e. -the  
persistence attributes, are adjusted, instead of focusing on the mean or the distribution.

~~As a first step in this method, QDM is applied separately on variables to adjust the empirical CDFs. This is followed by a  
385 multivariate bias adjustment to adjust the lag-0 and lag-1 auto and cross-correlation coefficients. This combination of univariate  
and multivariate bias-adjusting methods is applied on all time scales. For the multivariate adjustment, two models are used: a  
multivariate first-order autoregressive (AR(1)) model with constant parameters at the daily and yearly level, and a multivariate  
AR(1) model with periodic parameters (Salas, 1980) at the monthly and seasonal level. All steps are applied to the different  
types of data: historical observations of temperature, evaporation and precipitation (combined in the matrix  $\mathbf{X}^{\text{ho}}$ ), historical  
390 climate model simulations of the three variables (the matrix  $\mathbf{X}^{\text{hs}}$ ) and climate model projections of the three variables (the  
matrix  $\mathbf{X}^{\text{fs}}$ ), which have to be adjusted. All these datasets are of size  $T \times N$ , with  $T$  the number of time steps and  $N$  the  
number of variables.~~

The quantile-mapped future GCM time series for time step  $t$  is denoted as  $\mathbf{X}_t^{\text{fa}}$ . The standardised versions of this time series  
and of the observed time series are denoted as  $\check{\mathbf{X}}_t^{\text{fa}}$  and  $\check{\mathbf{X}}_t^{\text{ho}}$ , respectively. Using the standardised time (zero mean and unit  
395 variance) series, the multivariate AR(1) model with constant parameters (MAR) for the observed and GCM multivariate time

**Table 2.** Overview of the multivariate bias-adjusting methods

|                              | <u>MBCn</u>   | <u>MRQNBC</u>  | <u>R<sup>2</sup>D<sup>2</sup></u>                            | <u>dOTC</u>   |
|------------------------------|---|--|--|---|
| <u>Category</u>              | <u>Marginal/dependence</u>  | <u>Marginal/dependence</u>   | <u>Marginal/dependence</u>                                   | <u>All-in-one</u>                                   |
| <u>Temporal properties</u>   | <u>Shuffle based on observations</u>  | <u>Observed</u>  | <u>Shuffle based on observations</u>                         | <u>Future, adjusted</u>                             |
| <u>Dependence structure</u>  | <u>Future, adjusted based on observations</u>                                     | <u>Observed</u>  | <u>Observed</u>  | <u>Future, adjusted</u>                             |
| <u>Trend preservation</u>    | <u>Marginal properties by the application of QDM, dependence structure partly</u> | <u>Marginal properties by the application of QDM</u>                 | <u>Marginal properties by the application of QDM</u>         | <u>Marginal properties and dependence structure</u> |
| <u>Statistical technique</u> | <u>Iterative partial matrix recorrelation</u>                                     | <u>Autoregressive modeling</u>                                       | <u>Conditional resampling</u>                                | <u>Optimal transport</u>                            |
| <u>Timescale</u>             | <u>Daily adjustment by QDM + full time series shuffle</u>                         | <u>Combination of daily, monthly, seasonal and yearly adjustment</u> | <u>Daily adjustment by QDM + full time series resampling</u> | <u>Full time series</u>                             |

series can be expressed as (Mehrotra and Sharma, 2016):-

$$\check{X}_t^{\text{ho}} = \mathbf{C}\check{X}_{t-1}^{\text{ho}} + \mathbf{D}\epsilon_t$$

and

$$\check{X}_t^{\text{fa}} = \mathbf{E}\check{X}_{t-1}^{\text{fa}} + \mathbf{F}\epsilon_t,$$

400 with  $\mathbf{C}$  and  $\mathbf{D}$  the coefficient matrices of  $\check{X}_t^{\text{ho}}$ ,  $\mathbf{E}$  and  $\mathbf{F}$  the coefficient matrices of  $\check{X}_t^{\text{fa}}$  and  $\epsilon_t$  a white noise term. The coefficient matrices are calculated using the  $N \times N$  lag-0 and lag-1 cross-correlation matrices  $\mathbf{M}_0$  and  $\mathbf{M}_1$ . Using the standardised time

series, the elements of these matrices can be expressed as (Salas, 1980):-

$$\underline{m}_0^{i,j} = \frac{1}{T} \sum_{t=1}^T x_t^i x_t^j$$

$$\underline{m}_1^{i,j} = \frac{1}{T-1} \sum_{t=1}^T x_t^i x_{t+1}^j,$$

405 with  $i$  and  $j$  the column numbers of  $\check{\mathbf{X}}_t$ , referring to the variables whose correlation is calculated. This enables the calculation of  $\mathbf{C}$  and  $\mathbf{E}$  as (Matalas, 1967):-

$$\underline{\mathbf{C}} = \underline{\mathbf{M}}_1^{\text{ho}} \underline{\mathbf{M}}_0^{\text{ho}-1},$$

$$\underline{\mathbf{E}} = \underline{\mathbf{M}}_1^{\text{fa}} \underline{\mathbf{M}}_0^{\text{fa}-1},$$

and of  $\mathbf{D}$  and  $\mathbf{F}$  via-

$$410 \underline{\mathbf{D}} \underline{\mathbf{D}}^T = \underline{\mathbf{M}}_0^{\text{ho}} - \underline{\mathbf{M}}_1^{\text{ho}} \underline{\mathbf{M}}_0^{\text{ho}-1} \underline{\mathbf{M}}_1^{\text{ho}T}$$

$$\underline{\mathbf{F}} \underline{\mathbf{F}}^T = \underline{\mathbf{M}}_0^{\text{fa}} - \underline{\mathbf{M}}_1^{\text{fa}} \underline{\mathbf{M}}_0^{\text{fa}-1} \underline{\mathbf{M}}_1^{\text{fa}T},$$

which can be solved using the eigenvalues and eigenvectors of  $\underline{\mathbf{D}} \underline{\mathbf{D}}^T$  or  $\underline{\mathbf{F}} \underline{\mathbf{F}}^T$ :-

$$\underline{\mathbf{D}} = \underline{\mathbf{V}} \sqrt{\underline{\mathbf{S}}} \underline{\mathbf{V}}^T,$$

$$\underline{\mathbf{F}} = \underline{\mathbf{V}} \sqrt{\underline{\mathbf{S}}} \underline{\mathbf{V}}^T,$$

415 with  $\underline{\mathbf{V}}$  the matrix of eigenvectors and  $\underline{\mathbf{S}}$  a diagonal matrix with the corresponding eigenvalues-

The multivariate bias adjustment is then implemented by removing the lag-0 and lag-1 auto- and cross-correlations from the future time series  $\check{\mathbf{X}}_t^{\text{fa}}$  (the matrices  $\mathbf{E}$  and  $\mathbf{F}$ ) and applying the observed lag-0 and lag-1 auto- and cross-correlations ( $\mathbf{C}$  and  $\mathbf{D}$ ) to the future time series and thus creating a modified future time series,  $\check{\mathbf{X}}_t^{\prime \text{fa}}$ . These steps are applied by first rearranging and simplifying Eq. (??) for  $\epsilon_t$ :-

$$420 \underline{\epsilon}_t = \underline{\mathbf{F}}^{-1} \left( \check{\mathbf{X}}_t^{\text{fa}} - \underline{\mathbf{E}} \check{\mathbf{X}}_{t-1}^{\text{fa}} \right),$$

with  $\underline{\epsilon}_t$  now a standardised vector of  $N$  variables calculated by removing the lag-0 and lag-1 auto- and cross-correlations from the  $\check{\mathbf{X}}_t^{\text{fa}}$  time series. This vector is plugged into Eq. (??) along with the matrices  $\mathbf{C}$  and  $\mathbf{D}$  in which  $\check{\mathbf{X}}_t^{\text{fa}}$  is used instead of  $\check{\mathbf{X}}_t^{\text{ho}}$  to obtain the modified time series:-

$$\underline{\check{\mathbf{X}}_t^{\prime \text{fa}}} = \underline{\mathbf{C}} \underline{\check{\mathbf{X}}_{t-1}^{\text{fa}}} + \underline{\mathbf{D}} \underline{\mathbf{F}}^{-1} \left( \check{\mathbf{X}}_t^{\text{fa}} - \underline{\mathbf{E}} \check{\mathbf{X}}_{t-1}^{\text{fa}} \right),$$

425 which can be rearranged as:-

$$\underline{\check{\mathbf{X}}_t^{\prime \text{fa}}} = \underline{\mathbf{C}} \underline{\check{\mathbf{X}}_{t-1}^{\text{fa}}} + \underline{\mathbf{D}} \underline{\mathbf{F}}^{-1} \check{\mathbf{X}}_t^{\text{fa}} - \underline{\mathbf{D}} \underline{\mathbf{F}}^{-1} \underline{\mathbf{E}} \check{\mathbf{X}}_{t-1}^{\text{fa}}.$$

This model preserves the The biases are adjusted by replacing the modeled persistence attributes with observed persistence attributes. As a last step, destandardising results in the bias-adjusted time series  $\check{\mathbf{X}}_t^{\prime \text{fa}}$ .

When using the multivariate AR(1) with periodic parameters (PMAR), the parameters are derived separately for each period to allow for periodicity. In this case, the vectors  $\mathbf{X}_{t,\tau}^{\text{ho}}$  and  $\mathbf{X}_{t,\tau}^{\text{fa}}$  respectively represent the observed and quantile mapped-GCM time series. The subscript  $t$  refers to the year and the subscript  $\tau$  to a specific period in the year.

The elements of the periodic version of  $\mathbf{M}_0$  and  $\mathbf{M}_1$  can be calculated as (Salas, 1980):-

$$m_{0,\tau}^{i,j} = \frac{\sum_{t=1}^{T_\tau} (x_{t,\tau}^i - \bar{x}_\tau^i)(x_{t,\tau}^j - \bar{x}_\tau^j)}{T_\tau s_\tau^i s_\tau^j}$$

$$m_{1,\tau}^{i,j} = \frac{\sum_{t=1}^{T_\tau} (x_{t,\tau}^i - \bar{x}_\tau^i)(x_{t,\tau-1}^j - \bar{x}_{\tau-1}^j)}{T_\tau s_\tau^i s_{\tau-1}^j},$$

with  $T_\tau$  the number of time steps of the period  $\tau$ ,  $\bar{x}_\tau$  and  $\bar{x}_{\tau-1}$  the mean of periods  $\tau$  and  $\tau-1$  (for instance, if  $\tau$  is summer, than  $\tau-1$  is spring) and  $s_\tau$  and  $s_{\tau-1}$  the standard deviations of periods  $\tau$  and  $\tau-1$ . The correlation matrices are calculated in the same way as in the non-periodic steps. The only difference is that they are calculated for every period (e.g. separately for every season or month). For every time step in period  $\tau$ , the corresponding value can be adjusted as follows to preserve the observed persistence attributes:-

$$\check{\mathbf{X}}_{t,\tau}^{\text{fa}} = \mathbf{C}_\tau \check{\mathbf{X}}_{t,\tau-1}^{\text{fa}} + \mathbf{D}_\tau \mathbf{F}_\tau^{-1} \check{\mathbf{X}}_{t,\tau}^{\text{fa}} - \mathbf{D}_\tau \mathbf{F}_\tau^{-1} \mathbf{E}_\tau \check{\mathbf{X}}_{t,\tau-1}^{\text{fa}}.$$

The different time steps are combined with the nesting method proposed in Johnson and Sharma (2012) and Mehrotra and Sharma (2015). First, QDM (as described in Section 3.2.1) is applied at a daily level, followed by MAR. These adjusted time series are then aggregated and averaged to form a monthly time series, which is adjusted by QDM, standardised and adjusted by PMAR. Note that the standardisation of the aggregated time series does not imply that the variables of a period  $\tau$  of that time series have zero mean and unit variance. The results of the monthly adjustment are aggregated and averaged to form seasonal time series, which are also adjusted using QDM, standardised and adjusted by PMAR. As a last nesting step, the results are once more aggregated and averaged to build an annual time series, which is adjusted using QDM and MAR. The outcomes of all these steps are combined into a weighting factor that is used to modify the daily time series accordingly (Srikanthan and Pegram, 2009):-

$$\check{\mathbf{X}}_{t,j,s,i}^{\text{ffa}} = \left( \frac{\mathbf{Y}_{j,s,i}^{\text{fa}}}{\mathbf{Y}_{j,s,i}^{\text{fa}}} \right) \left( \frac{\mathbf{Z}_{s,i}^{\text{fa}}}{\mathbf{Z}_{s,i}^{\text{fa}}} \right) \left( \frac{\mathbf{A}_i^{\text{fa}}}{\mathbf{A}_i^{\text{fa}}} \right) \mathbf{X}_{t,j,s,i}^{\text{fa}},$$

with  $t$  the day,  $j$  the month,  $s$  the season,  $i$  the year,  $\mathbf{Y}_{j,s,i}^{\text{fa}}$  the monthly adjusted value,  $\mathbf{Y}_{j,s,i}^{\text{fa}}$ , on the basis of autoregressive expressions. Besides replacing the simulated temporal properties with the aggregated-averaged monthly value,  $\mathbf{Z}_{s,i}^{\text{fa}}$  the seasonal adjusted value,  $\mathbf{Z}_{s,i}^{\text{fa}}$  observed ones, this implies that the simulated dependence structure is also replaced by the observed structure. As QDM is applied on each timescale, the aggregated-averaged seasonal value,  $\mathbf{A}_i^{\text{fa}}$  the adjusted yearly value and  $\mathbf{A}_i^{\text{fa}}$  the aggregated-averaged yearly value. The full procedure is summarised in Algorithm ??, marginal properties are preserved.

MRQNBC Daily historical observations  $\mathbf{X}^{\text{ho}}$  Daily historical simulations  $\mathbf{X}^{\text{hs}}$  Daily future simulations  $\mathbf{X}^{\text{fs}}$  Adjusted future simulations  $\check{\mathbf{X}}^{\text{ffa}}$  Apply QDM to calculate  $\mathbf{X}^{\text{fa}}$  Standardise  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{fa}}$  Calculate matrices  $\mathbf{C}$  and  $\mathbf{D}$  of  $\check{\mathbf{X}}^{\text{ho}}$  and  $\mathbf{E}$  and  $\mathbf{F}$



~~of  $\tilde{X}^{fa}$ . Apply the persistence adjustment to calculate  $\tilde{X}^{fa}$ . Destandardise  $\tilde{X}^{fa}$ . Aggregate  $X^{ho}$  and  $X^{fa}$  to the higher timescale. Calculate weighting factors for all timescales except the daily timescale. Calculate the final adjusted daily value  $X^{fa}$ .~~

460 ~~The~~ After adjusting all timescales, the final daily result is calculated by weighing all timescales. However, as the nesting method cannot fully remove biases at all time scales, ~~thus~~ Mehrotra and Sharma (2016) suggested to repeat the **complete entire** procedure multiple times. ~~However~~ Yet, in our case ~~this seemed to exacerbate the results, so the method was run only once~~ multiple repetitions exacerbated the results. Non-realistic outliers created by the first repetition influenced the subsequent repetitions, creating even more non-realistic values. This was most clearly seen for precipitation. As a bounded variable, precipitation is most sensitive for non-realistic values. Nonetheless, running the method just once yielded satisfactory results.  
465 A full mathematical description of the method can be found in Mehrotra and Sharma (2016).

### 3.3.2 Multivariate Bias Correction in $n$ dimensions

In 2018, Cannon (2018) proposed the ‘Multivariate Bias correction in  $n$  dimensions’ (MBCn) method as a flexible multivariate bias-adjusting method. The method’s flexibility has attracted some attention, ~~as~~ and it has already been used in multiple studies (Räty et al., 2018; Zscheischler et al., 2019; Meyer et al., 2019; François et al., 2020). This method consists of three steps.  
470 First, the multivariate data are rotated using a randomly generated orthogonal rotation matrix, adjusted with the additive form of QDM, and rotated back until the calibration period model simulations converge to the observations. This convergence is verified on the basis of the energy distance (Rizzo and Székely, 2016). Second, the validation period simulations are adjusted using QDM, as this method preserves the simulated trends. As the last step, these adjusted time series are shuffled using the Schaake Shuffle (Clark et al., 2004) based on the rank order of the rotated dataset. ~~A thorough mathematical explanation full~~  
475 mathematical description of the method can be found in Cannon (2018) ~~and Van de Velde et al. (2020), and the implementation is~~.

### 3.3.3 Rank Resampling for Distributions and Dependences

One of the most recent methods studied in this paper is the ‘Rank Resampling for Distributions and Dependences’ ( $R^2D^2$ ) method, which was designed by Vrac (2018) as an improvement of the older EC-BC method (Vrac and Friederichs, 2015).  
480 Recently,  $R^2D^2$  was further extended for better multisite and temporal representation by Vrac and Thao (2020b) ( $R^2D^2$  v2.0). This method is a marginal/dependence multivariate bias-adjusting method, which adjusts the simulated climate dependence by resampling from the observed dependence. The resampling is applied through the search for an analogue for the ranks of a simulated reference dimension in the observed time series, which makes this an application of the analogue principle (Lorenz, 1969; Zorita and Von Storch, 1999) in bias adjustment. A detailed mathematical description can be found in Vrac (2018)  
485 and Vrac and Thao (2020b).

In the present application of  $R^2D^2$ , QDM was used as the univariate bias-adjusting method to ensure consistency with the other multivariate bias-adjusting methods. This ensures the preservation of the changes in the same as in the latter. For the sake of clarity, marginal distribution, besides the preservation of some temporal properties, which is inherent to the method

is summarised in Algorithm ??-method. Each variable (precipitation, evaporation and temperature) was in turn used as the reference dimension. No other information was included, as the present study was limited to one grid cell.

MBCn Historical observations  $\mathbf{X}^{\text{ho}}$  Historical simulations  $\mathbf{X}^{\text{hs}}$  Future simulations  $\mathbf{X}^{\text{fs}}$  Adjusted future simulations  $\mathbf{X}^{\text{fa}}$   
 Initialisation: tolerance  $\epsilon$  and initial energy distance difference  $\Delta D_0$  Randomly generate a rotation matrix  $\mathbf{R}$  Rotate  $\mathbf{X}^{\text{ho}}$ ,  $\mathbf{X}^{\text{hs}}$  and  $\mathbf{X}^{\text{fs}}$  Apply the additive form of QDM Rotate  $\mathbf{X}^{\text{ho}}$ ,  $\mathbf{X}^{\text{hs}}$  and  $\mathbf{X}^{\text{fs}}$  back Calculate the energy distance  $D$  between  $\mathbf{X}^{\text{hs}}$  and  $\mathbf{X}^{\text{ho}}$  Calculate the decrease in energy distance  $\Delta D$  Apply QDM to the original inputs to calculate  $\mathbf{X}^{\text{fa}}$  Apply the Schaake  
 Shuffle based on the rotated future simulations to calculate  $\mathbf{X}^{\text{fa}}$

### 3.3.4 Dynamical Optimal Transport Correction

Recently, Robin et al. (2019) indicated that the notion of a transfer function in quantile mapping can be generalised to the theory of optimal transport. Optimal transport is a way to measure the dissimilarity between two probability distributions and to use this as a means for transforming the distributions in the most optimal way (Villani, 2008; Peyré and Cuturi, 2019).

Optimal transport was used by Robin et al. (2019) to adjust the bias of a multivariate data set in the ‘dynamical Optimal Transport Correction’ method (dOTC), which extends the ‘CDF-transform’ (CDF-t) bias-adjusting method (Michelangeli et al., 2009) to the multivariate case. dOTC calculates the optimal transport plans from  $\mathbf{X}^{\text{ho}}$  to  $\mathbf{X}^{\text{hs}}$  (the bias between the model and the simulations) and from  $\mathbf{X}^{\text{hs}}$  to  $\mathbf{X}^{\text{fs}}$  (the evolution of the model). The combination of both optimal transport plans allows for bias adjustment while preserving the trend of the model.

The different transport plans and transformations used in dOTC. Based on Robin et al. (2019).

Optimal transport is applied on the basis of an optimal transport plan. The optimal plan between  $\mathbf{X}^{\text{ho}}$  simulated changes in both marginal properties and  $\mathbf{X}^{\text{hs}}$  is denoted as  $\gamma$ . The second optimal plan between  $\mathbf{X}^{\text{hs}}$  and  $\mathbf{X}^{\text{fs}}$  is denoted as  $\phi$ . The goal is to transform  $\phi$  according to  $\gamma$ , defining a new plan  $\tilde{\phi}$ . This optimal plan estimates  $\mathbf{X}^{\text{fa}} = \tilde{\phi}(\mathbf{X}^{\text{ho}})$ . Finally,  $\mathbf{X}^{\text{fs}}$  is adjusted with respect to  $\mathbf{X}^{\text{fa}}$ , creating the final adjusted  $\mathbf{X}^{\text{fa}}$ . These steps are summarised in Fig. ??.

For the definition of the optimal plan, the ‘Optimal Transport Correction’ (OTC) (Robin et al., 2019) is used. First, the empirical distributions  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$  and  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$  have to be calculated. To achieve this, the subspace of  $\mathbb{R}^N$  that contains the data is partitioned into regularly spaced cells, generically denoted  $\mathbf{c}_i^*$ , with  $N$  the number of variables of  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{hs}}$ . Using this notation,  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$  and  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$  can be estimated using the relative frequencies  $p_{\mathbf{X}^{\text{ho}}}$  and  $p_{\mathbf{X}^{\text{hs}}}$  as:

$$p_{\mathbf{X}^{\text{ho}},i} = \frac{1}{m} \sum_{k=1}^m \mathbf{1}(\mathbf{X}_k^{\text{ho}} \in \mathbf{c}_i^*),$$

$$p_{\mathbf{X}^{\text{hs}},i} = \frac{1}{n} \sum_{l=1}^n \mathbf{1}(\mathbf{X}_l^{\text{hs}} \in \mathbf{c}_i^*),$$

with  $\mathbf{1}$  the indicator function and  $m$  and  $n$  the total number of time steps of respectively  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{hs}}$ . Thus, the distributions are essentially estimated by counting the number of observations of each time series within each cell. The optimal plan  $\gamma$  between  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{hs}}$  can be estimated as:

$$\hat{\gamma} = \sum_{i,j=1}^{I,J} \gamma_{i,j}.$$

520 The coefficients  $\gamma_{i,j}$  are the probabilities to transform an observation of  $\mathbf{X}^{\text{ho}}$  in cell  $\mathbf{e}_i^*$  into an observation of  $\mathbf{X}^{\text{hs}}$  in cell  $\mathbf{e}_j^*$  and  $I$  and  $J$  are the total number of cells containing observations of respectively  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{hs}}$ . Note that from here on,  $\mathbf{e}_i^*$  and  $\mathbf{e}_j^*$  denote only those cells containing observations of respectively  $\mathbf{X}^{\text{ho}}$  and  $\mathbf{X}^{\text{hs}}$  and that ‘observation’ is used generically for both observed and simulated time series. The coefficients are unknown, but obey the marginal properties :-

$$\sum_{j=1}^J \gamma_{i,j} = p_{\mathbf{X}^{\text{ho}},i}$$

525 and  $\sum_{i=1}^I \gamma_{i,j} = p_{\mathbf{X}^{\text{hs}},j}.$

Central to the optimal transport theorem is the cost function  $C$  (Villani, 2008), which can here be approximated by-

$$\hat{C}(\hat{\gamma}) = \sum_{i,j=1}^{I,J} \|\mathbf{c}_i - \mathbf{c}_j\|^2 \gamma_{i,j},$$

with  $\|\cdot\|$  the Euclidean norm and  $\mathbf{c}_i$  and  $\mathbf{c}_j$  the centres of the cells defined above. Finding  $\gamma_{i,j}$  comes down to solving the problem defined by the constraints of Eqs. (??) and minimisation of Eq. (??). Here, Sinkhorn’s algorithm (Cuturi, 2013) is used to find the solution to this problem and thus the optimal plan  $\gamma$ . Using this optimal plan, a vector of probabilities with length  $J$  can be defined for each cell  $\mathbf{e}_i^*$  of  $\mathbf{X}^{\text{ho}}$ , each element  $j$  corresponding to the probability that an observation in that cell  $\mathbf{e}_i^*$  will be transformed into an observation in cell  $\mathbf{e}_j^*$  of  $\mathbf{X}^{\text{hs}}$ . In the transformation, the vectors of probabilities are used to introduce stochasticity, by sampling from these vectors the element  $j$  corresponding to cell  $\mathbf{e}_j^*$ . The stochastic transformation of an observation of  $\mathbf{X}^{\text{ho}}$  into an observation of  $\mathbf{X}^{\text{hs}}$  can be repeated to create an ensemble of results. This ensemble accounts for random weather effects and can thus be considered to be more similar to the true range of observations-

535

The optimal plan  $\phi$  can be calculated analogously. This optimal plan  $\phi$  transforms an observation of  $\mathbf{X}^{\text{hs}}$  in cell  $\mathbf{e}_j^*$  into an observation of  $\mathbf{X}^{\text{fs}}$  in cell  $\mathbf{e}_k^*$ , with  $\mathbf{e}_k^*$  defined analogously to  $\mathbf{e}_i^*$  and  $\mathbf{e}_j^*$ . What distinguishes dOTC from OTC is the next phase, in which  $\phi$  is transformed according to  $\gamma$ , resulting in  $\tilde{\phi}$ . This is conducted in three steps, the first being is the transformation of  $\phi$  into a vector. The vector  $\mathbf{v}_{jk} := \mathbf{e}_k - \mathbf{e}_j$  represents the climatic trend from an observation of  $\mathbf{X}^{\text{hs}}$  in cell  $\mathbf{e}_j^*$  to an observation of  $\mathbf{X}^{\text{fs}}$  in cell  $\mathbf{e}_k^*$ . The second step is the transfer according to  $\gamma$ . The result  $\tilde{\phi}$  can be defined by translating the observations of  $\mathbf{X}^{\text{ho}}$  along their respective vectors  $\mathbf{v}_{jk}$ : an observation of  $\mathbf{X}^{\text{fa}}$  is then given by  $\mathbf{X}_t^{\text{ho}} + \mathbf{v}_{jk}$ , with  $\mathbf{X}_t^{\text{ho}}$  the observation of  $\mathbf{X}^{\text{ho}}$  at time step  $t$ . However, the translation of  $\mathbf{X}_t^{\text{ho}}$  along vector  $\mathbf{v}_{jk}$  does not always define an optimal transport plan: the vector has to be adapted to  $\mathbf{X}^{\text{ho}}$ , which is the third step. In this step, a matrix factor  $\mathbf{D}$  is introduced, which rescales the vector  $\mathbf{v}_{jk}$ . This

540

545 resealing is actually the replacement of the scale of  $\mathbf{X}^{\text{hs}}$  by that of  $\mathbf{X}^{\text{ho}}$ . A Cholesky decomposition of the covariance matrix has been proposed for this resealing (Bárdossy and Pegram, 2012; Cannon, 2016). Denoting the covariance matrix as  $\Sigma$ , and the Cholesky decomposition as  $\text{Cho}(\Sigma)$ , Robin et al. (2019) proposed to multiply  $\mathbf{v}_{jk}$  by the following matrix:

$$\mathbf{D} := \text{Cho}(\Sigma_{\mathbf{X}^{\text{ho}}}) \cdot \text{Cho}(\Sigma_{\mathbf{X}^{\text{hs}}})^{-1}.$$

550 Robin et al. (2019) remark that the Cholesky decomposition only exists if  $\Sigma$  is symmetric and positive-definite. Some covariance matrices, such as those of highly correlated random variables, do not have this property.  $\Sigma$  must then be slightly perturbed to be positive-definite (Higham, 1988; Knol and ten Berge, 1989). It is also possible for the Cholesky decomposition to be poorly estimated if the available data are too small compared to the dimension. In that case, it is suggested to replace the matrix  $\mathbf{D}$  by the diagonal matrix of the standard deviation:  $\mathbf{D} = \text{diag}(\sigma_{\mathbf{X}^{\text{ho}}}\sigma_{\mathbf{X}^{\text{hs}}}^{-1})$ .

555 An observation of  $\mathbf{X}^{\text{fa}}$  is then given by  $\mathbf{X}_t^{\text{ho}} + \mathbf{D} \cdot \mathbf{v}_{jk}$ . To finalize, the empirical distribution  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fa}}}$  can be calculated. Using this distribution, OTC can be applied to  $\mathbf{X}^{\text{hs}}$  and  $\mathbf{X}^{\text{fa}}$  to generate  $\mathbf{X}^{\text{fa}}$ . A more elaborate mathematical explanation can be [the dependence structure](#). A full mathematical description of the method can be found in Robin et al. (2019), a summary is given in Algorithm ??.

560 ~~dOTC Historical observations  $\mathbf{X}^{\text{ho}}$  Historical simulations  $\mathbf{X}^{\text{hs}}$  Future simulations  $\mathbf{X}^{\text{fs}}$  Adjusted future simulations  $\mathbf{X}^{\text{fa}}$  Calculate the empirical distributions  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$ ,  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$  and  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fs}}}$  Calculate the optimal plan  $\gamma$  between  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$  and  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$  Calculate the optimal plan  $\phi$  between  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fs}}}$  and  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$  Calculate the Cholesky factor  $\mathbf{D}$  Find cell  $e_i$  containing  $\mathbf{X}_t^{\text{ho}}$  Construct the vector of probabilities  $\hat{\gamma}_{\mathbf{X}_t^{\text{ho}}} = (\gamma_{i,1}, \dots, \gamma_{i,J}) / p_{\mathbf{X}^{\text{ho}},i}$  Sample  $j \in \{1, \dots, J\}$  according to the probability vector  $\hat{\gamma}_{\mathbf{X}_t^{\text{ho}}}$  Construct the vector of probabilities  $\hat{\phi}_{\mathbf{X}_t^{\text{hs}}} = (\phi_{j,1}, \dots, \phi_{j,K}) / p_{\mathbf{X}^{\text{hs}},j}$  Sample  $k \in \{1, \dots, K\}$  according to the probability vector  $\hat{\phi}_{\mathbf{X}_t^{\text{hs}}}$  Calculate the vector  $\mathbf{v}_{jk}$  Calculate  $\mathbf{X}_t^{\text{fa}} = \mathbf{X}_t^{\text{ho}} + \mathbf{D} \cdot \mathbf{v}_{jk}$  Calculate the empirical distribution  $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fa}}}$  Apply OTC between  $\mathbf{X}^{\text{fa}}$  and  $\mathbf{X}^{\text{fs}}$  to generate  $\mathbf{X}^{\text{fa}}$~~

### 3.4 Experimental design

565 Prior to all intensity-bias-adjusting methods, the thresholding occurrence-adjusting method was applied. [As in Van de Velde et al. \(2020\)](#) [In the intensity-bias-adjustment step](#), a balance was sought between randomness and computational power for the calculation of the intensity-bias-adjusting methods. Methods with randomised steps were repeated. As such, 10 calculations were made for dOTC. The resulting values of each index were averaged for further comparison. Biases on the indices were always calculated as raw or adjusted simulations minus observations, indicating a positive bias if the raw or adjusted simulations are larger than  
570 the observations and a negative bias if the simulations are smaller.

## 4 Results

In this section, the results will be shown first for the R index calculations for bias change, and then for the validation indices. For the validation indices, first the indices based on the adjusted variables are discussed, followed by an elaboration on the indices based on the derived variables. As the effect on discharge is the overarching goal of this paper and the discharge indices are

575 affected by all other indices, those will be discussed last. ~~All observed values and biases of both raw and adjusted simulations are presented in Table ??.~~

#### 4.1 Bias change

The ~~R index values for the variable averages, standard deviations and all indices are given in Table ??.~~ The results results for the R index vary considerably depending on the ~~variable and/or index: for P, the bias can be considered almost stationary: only the~~ 99.5th percentile has an R index above one. In contrast, for E, the season: bias nonstationarity (R index values > 1) is present for all variables, but the extent varies. For precipitation, bias nonstationarity is most clear in winter and summer for the highest percentiles ( $P_{99}$  and  $P_{99.5}$ ). For temperature, winter, spring and summer all show some high R index values ~~are above 1 for,~~ but while winter has high R index values for all percentiles, the nonstationarity is restricted to the lower to middle percentiles ( $T_5$ ,  $T_{25}$ ,  $T_{50}$  and  $T_{75}$ ) for spring and the lower percentiles ( $T_5$  and  $T_{25}$ ) for summer. This is reflected in the middle percentiles ~~and for the standard deviation, indicating some major changes in parts of the distribution, and consequently, the bias. For T, the mean and standard deviation: both are nonstationary for winter, whereas only the mean is nonstationary for spring and neither mean nor standard deviation is nonstationary for summer. In autumn, the behavior is less clear: two percentiles ( $P_{50}$  and  $P_{95}$ ) have an R index value of 2, but unlike the other seasons, there is no apparent pattern as these values are far apart. However, the lower extremes are clearly influenced, although the bias on the higher extremes does not change. The different effects on~~ the variables are linked with the effect on the (cross-)correlations. For example, the lag-1 cross-correlation between P and E standard deviation has an R index value higher than 1 for autumn temperatures, indicating that some bias nonstationarity could be present. For evaporation, spring has the clearest bias nonstationarity: almost all percentiles have an R index value higher than 1. For the other seasons, the nonstationarity is less striking, although present. For winter and autumn,  $E_{75}$  has an R index value of only 0.19, whereas the ~~1 or higher and a clearly nonstationary standard deviation, while in summer,  $E_{25}$  and  $E_{50}$~~  have an R index value for the cross-correlation between E and T is 1.20. Although the R index values are low for P, this does not imply that ~~higher than 1, although neither mean nor standard deviation is clearly nonstationary. For occurrence the bias nonstationarity seems limited: only in spring and autumn, the R index values for the precipitation occurrence indices are low. With an R index value of 1.44, the auto-correlation bias clearly changes between both periods. However, this is not reflected by the other precipitation occurrence indices, which all but one have R index values lower than one value for precipitation lag-1~~ autocorrelation is higher than 1. For correlation, the bias nonstationarity is also limited, although some of the correlations of evaporation and either temperature or precipitation have an R index value higher than 1, but this depends on the season (crosscorr $_{E,T,0}$  and crosscorr $_{E,T,1}$  in spring, crosscorr $_{E,T,1}$  in winter, corr $_{E,T}$  in summer and corr $_{P,E}$  in autumn).

R index values for 1970-1989 as historical period and 1998-2017 as future period **Indices R index Indices R index Indices R index Indices**

605  ~~$P_5$  NaN  $T_5$  2  $E_5$  0.03  $P_{lag1}$  1.44  $P_{25}$  0  $T_{25}$  2  $E_{25}$  0.47  $P_{P00}$  0.09  $P_{50}$  0.10  $T_{50}$  2  $E_{50}$  1.47  $P_{P10}$  0.41  $P_{75}$  0.13  $T_{75}$  0.87  $E_{75}$  2  $N_{dry}$  0.29  $P_{90}$  0.19  $T_{90}$  0.31  $E_{90}$  1.14 corr $_{E,T}$  0.75  $P_{95}$  0.17  $T_{95}$  0.07  $E_{95}$  1 corr $_{P,E}$  0.20  $P_{99}$  0.58  $T_{99}$  0.19  $E_{99}$  0.47 corr $_{P,T}$  0  $P_{99.5}$  1.02  $T_{99.5}$  0.08  $E_{99.5}$  0.20 crosseorr $_{E,T,0}$  2  $P_{mean}$  0.18  $T_{mean}$  2  $E_{mean}$  1.06 crosseorr $_{E,T,1}$  0.90  $P_{std}$  0.72  $T_{std}$  0.50  $E_{std}$  2 crosseorr $_{P,E,0}$  0.31 crosseorr $_{P,E,1}$  0.13 crosseorr $_{P,T,0}$  0.10 crosseorr $_{P,T,1}$  0.09~~

Many of the R index values ~~presented in Table ?? thus~~ indicate that the bias changes between the two periods considered here (1970-1989 versus 1998-2017) might already be large enough to have an effect on the bias adjustment. As these periods  
610 are only separated by 10 years, this is an important indicator for the bias adjustment of late 21st century data, just as Chen et al. (2015) mentioned. ~~However, it does not suffice to calculate just a few of these R index values.~~ The results vary substantially among ~~variables and even for the percentiles of a variable under consideration: while the 5th T percentile has an R index value of 2, the value for the 95th percentile is only 0.07.~~ This seasons, variables and distributions of the variables. Although this could give an indication of ~~why the methods perform more poorly for the reason for poor performance for~~ some of these  
615 indices. ~~However, purely based on these results,~~ it is impossible to ~~say state~~ exactly what causes the bias nonstationarities purely based on these results. Possible causes could be that recent trends such as those in precipitation extremes (Papalexiou and Montanari, 2019) are poorly captured by the models, that limiting mechanisms such as soil moisture depletion (Bellprat et al., 2013) are poorly modelled or that natural variability ~~influences (Addor and Fischer, 2015) (Addor and Fischer, 2015)~~ influences the biases. However, discussing this in depth is out of the scope of the present study and deserves a separate study.  
620 In what follows, we will focus on the performance of the bias-adjusting methods and whether or not there is a link with these nonstationarities.

## 4.2 Precipitation amount

~~Figure ?? presents the  $RB_O$  and  $RB_{MB}$  values for the highest P percentiles. None of the residual bias values of the lower percentiles can be plotted as either the observations are 0 mm ( $P_5$  and  $P_{25}$ ) or the  $RB_O$  values are lower than zero ( $P_{50}$ ). The~~  
625 ~~percentiles could also have been plotted for wet days only (e.g. days with P higher than 0.1 mm/day), but as some methods change the number of dry days after the initial thresholding step, the dry days are also included in the calculation of the indices. This influences the  $RB_O$  and  $RB_{MB}$  values: they are generally higher when the dry days are not included. The Perkins Skill Score (PSS) for precipitation (Table 3) indicates that the PDFs of the observations and adjusted simulations agree rather well. These scores are very similar in the calibration and validation period. Only QDM and mQDM perform worse in every season,~~  
630 ~~whereas the change performance of the multivariate methods depends on the season. For dOTC, the result is better in the validation period than in the calibration period.~~

~~The  $RB_O$  and  $RB_{MB}$  values depict a very similar performance for QDM, mQDM and MBCn, but a different performance for MRQNBC and dOTC. The similar performance of the former three is unsurprising, as their adjustments of P are all very similar. The only difference between QDM and MBCn versus mQDM is the time series to which the adjustment was applied,~~  
635 ~~as the latter is based on the historical time series. QDM, mQDM and MBCn are consistently the best methods out of the five tested here, with the~~

~~The good performance for the full PDF contrasts with the bias adjustment of the extreme values. Figure 1 presents the  $RB_{MB}$  values for  $P_{75}$ ,  $P_{90}$ ,  $P_{95}$  and  $P_{99}$  all below 0.5 and the  $RB_O$  values also below 0.5. The performances of MRQNBC and dOTC are worse, but not poor either:  $P_{75}$ ,  $P_{90}$ ,  $P_{95}$  and  $P_{99}$  all have  $RB_O$  and  $RB_{MB}$  values lower than 1, but only for dOTC the majority~~  
640 ~~of them (P for the highest P percentiles in the validation period. The lower percentiles ( $P_{75}$ ,  $P_{90}$ ,  $P_{95}$  and  $P_{99}$  to P~~

~~$RB_{MB}$  versus  $RB_O$  for the precipitation indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.~~

**Table 3.** PSS values for precipitation in the calibration (Cal) and validation (Val) periods (%).

|                               | Winter |      | Spring |      | Summer |      | Autumn |      |
|-------------------------------|--------|------|--------|------|--------|------|--------|------|
|                               | Cal    | Val  | Cal    | Val  | Cal    | Val  | Cal    | Val  |
| QDM                           | 96.5   | 94.1 | 97.2   | 95.0 | 97.9   | 96.7 | 96.8   | 94.3 |
| mQDM                          | 100.0  | 92.1 | 100.0  | 93.4 | 100.0  | 96.5 | 100.0  | 94.8 |
| MBCn                          | 94.1   | 89.9 | 95.6   | 95.5 | 92.6   | 97.3 | 97.2   | 96.8 |
| MRQNBC                        | 93.2   | 92.8 | 84.7   | 81.0 | 96.1   | 95.6 | 93.5   | 91.0 |
| dOTC                          | 64.6   | 73.8 | 66.6   | 70.0 | 62.2   | 91.6 | 62.8   | 72.0 |
| R <sup>2</sup> D <sup>2</sup> | 93.3   | 90.1 | 94.0   | 93.1 | 93.2   | 92.6 | 94.1   | 93.8 |

A surprising result for P is the high  $RB_{MB}$  value for 50 are adjusted very well by all methods, but the performance of the methods for the higher percentiles differs considerably between the seasons. For winter (blue) and summer (yellow), only  $P_{99.5}$  for MRQNBC. This percentile is too biased for all other methods to be plotted: the result is possibly influenced by the need for the observed values to be extrapolated if the simulated value lies outside the range of observed values (see e.g. Li et al. (2010)). This extrapolation is not implemented in this version of QDM and dOTC, which clearly hampers the ability of all methods to correctly adjust this percentile. Yet, as the  $P_{75}$  and  $P_{90}$  can be plotted in the validation period, whereas for spring and autumn all percentiles, from  $P_{75}$  to  $P_{99.5}$  can be plotted for all methods. The poor adjustment of the high percentiles in winter and summer could be caused by bias nonstationarity: the R index values for these percentiles are higher than 1, in contrast with the low and well-adjusted higher percentiles for spring and autumn precipitation. However, although  $P_{95}$  has an R index value was close to lower than 1 for both winter and summer, it is possible that bias nonstationarity also slightly influences the performance. For MRQNBC, however, the combination of QDM with the focus on correlation seemingly improves the performance of this percentile. As heavy precipitation values are clustered in time, the performance of the respective indices might be improved by the correlation. The good representation of heavy precipitation values in the MRQNBC-adjusted time series is also shown by  $P_{99}$ , for which MRQNBC has an  $RB_{MB}$  value of 0.59, the best of all methods.

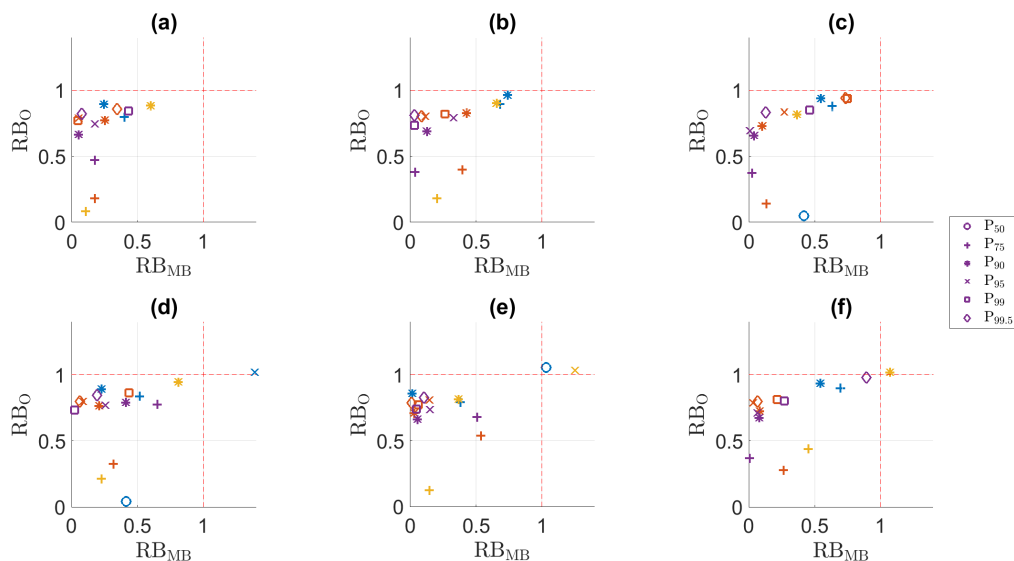
### 4.3 Temperature

For the temperature adjustment, the  $RB$  is poorly adjusted. This illustrates that the R index gives an indication of the nonstationarity, but also hides information on the size of the biases. For summer, the bias for  $P_{95}$  changes from 5.09 mm in the calibration period to 1.89 mm in the validation period, a change of over 3 mm. For winter, the bias changes from 1.44 mm in the calibration period to 0.52 mm in the validation period, a change of almost 1 mm. Yet, these differences have a very similar R index value. A comparison with the  $RB_{OMB}$  and  $RB_{MBO}$  values indicate that all methods result in a performance better than the raw climate simulations, except for MRQNBC of the calibration period (Fig. ??). In contrast to all other methods, only the residual bias values of  $T_{90}$  of MRQNBC are within the area delineated by the 1-1 lines. For all other indices, the bias is worse than in the climate simulations, with absolute biases up to 7°C. The results for MRQNBC are interesting, as T is the best understood

665 variable (Shepherd, 2014) and should thus not be hard to adjust. One line of reasoning could be the implementation of model trend preservation. While trend preservation is a prominent aspect in all other methods, the persistence preservation is at least as important in MRQNBC. The trade-off between both aspects of the bias adjustment thus seems to influence the result of MRQNBC, while the other methods can more easily adapt to and adjust the simulated T2) illustrates that all methods perform well for every season, indicating that the nonstationarity could be a cause of the diverging performances in the validation period between the winter/summer and spring/autumn pairs. However, this nonstationarity is not apparent from the PSS, as it only occurs in the tail of the distribution. This also follows from the R index values for the mean and standard deviation in winter and summer. Only for standard deviation, the R index value indicates nonstationarity in winter and summer: the values are respectively 1.79 and 1.56. Thus, the nonstationarity of the extremes and the standard deviation seem to be linked.

670

When comparing the indices for the other methods, the results are rather similar. They all have RB

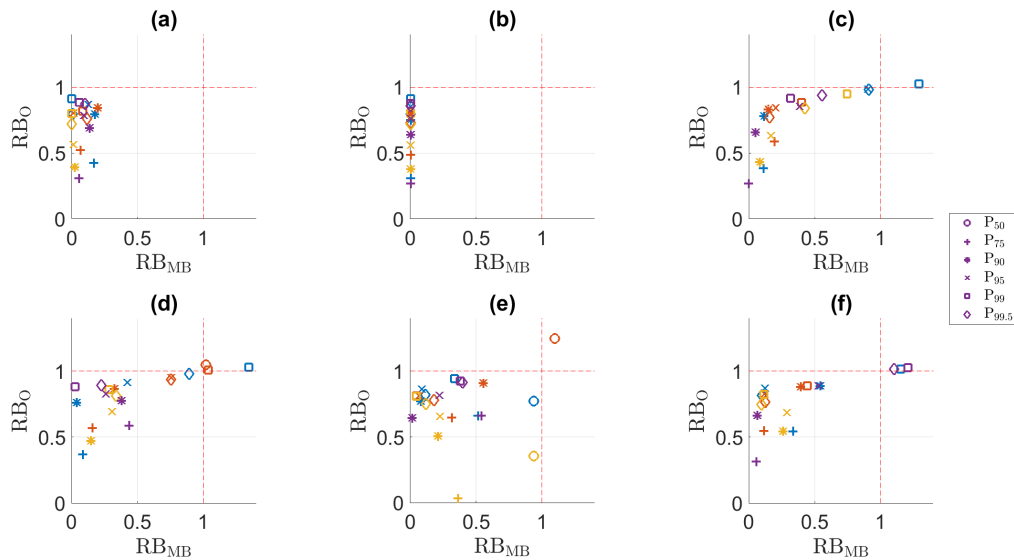


**Figure 1.**  $RB_{MB}$  versus  $RB_0$  for the precipitation in the validation period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.

675 The methods seem to perform rather similarly in every season. Although the  $RB_{eMB}$  values close to 1 vary, indicating that for some methods the bias is removed to a larger extent, the bias difference is small in comparison to the absolute T-values. Besides that, for every method, the lower T-percentiles have  $RB_{MB}$  values that are too high to be plotted, indicating that relative to the observations, the influence of the difference in removed bias is low. However, despite their similar behaviour, there is a difference that should be acknowledged. For example, on a yearly basis, the mean number of heavy precipitation days (R10, one of the ETCCDI indices (Zhang et al., 2011)) is well presented by all adjusted simulations (Fig. 3), but the methods show some notable differences. The highest percentiles have the lowest  $RB_{MB}$  values for QDM and MBCn, which have the same percentiles by construction, but this differs for the other methods. For example,  $T_{99.5}$  has an  $RB_{MB}$  value of 0.09 for QDM

680





**Figure 2.**  $RB_{MB}$  versus  $RB_0$  for the precipitation indices in the calibration period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.

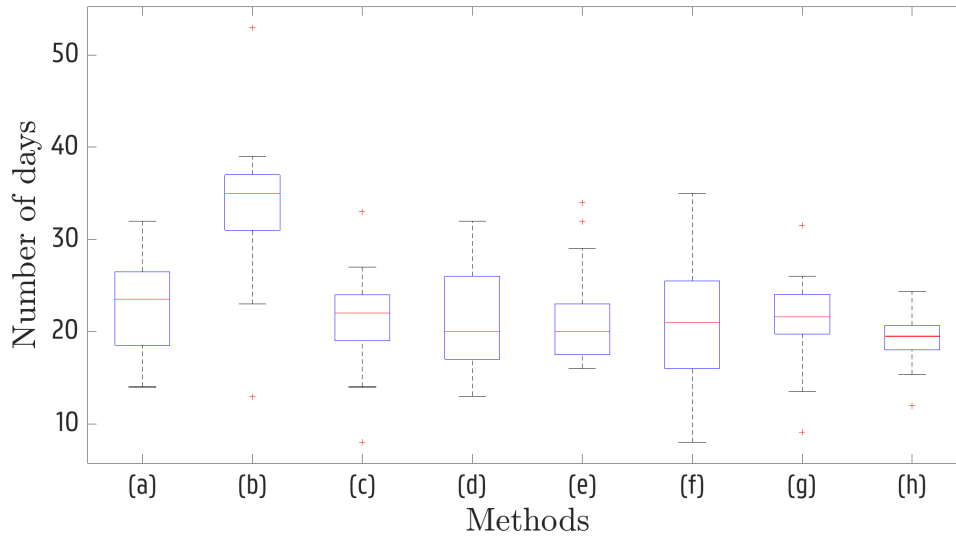
and MBCn, but only 0.43 for dOTC and 0.77 for mQDM. On the other hand, when considering all plotted percentiles, dOTC generally performs best. The highest  $RB_{MB}$  value for dOTC is 0.52 ( $T_{75}$ ), whereas 0.65 ( $T_{75}$ ) is the highest value for QDM and MBCn and 0.77 ( $T_{99.5}$ ) for mQDM. Although broadly similar, the indices for QDM and mQDM display some interesting differences. Whereas for QDM  $T_{95}$ ,  $T_{99}$  and  $T_{99.5}$  have the lowest  $RB_{MB}$  values,  $T_{75}$  has the lowest value for mQDM. In contrast, QDM has the highest  $RB_{MB}$  value for  $T_{75}$ . This might imply that for the highest T values, it is better to follow the simulations, while for slightly lower values, it is better to only use the climate change signal. Yet, QDM has the best  $RB_{MB}$  values and might thus be preferable. yearly variance clearly depends on the method: MRQNBC overestimates the variance, whereas the other methods slightly underestimate it.

For the lowest T values, all methods seem unable to handle the change in bias (as seen in Table ??): the  $RB_{MB}$  values are all higher than 1. This poor performance, combined with the high values for  $RB_0$ , might imply that it is better not to adjust T and work with the raw climate simulations. However, for the extreme T values-

### 4.3 Temperature

Table 4 displays the PSS values for temperature. It can be seen that the univariate bias-adjusting methods have higher values than the multivariate methods for all seasons. Among the multivariate methods, the performance also varies: dOTC performs best, the absolute biases can be more than 1°C. Thus, depending on the research goal and the R index value, it might be important to consider whether or not T should be adjusted, whereas the performance for the other multivariate bias-adjusting methods depends strongly on the season. However, the multivariate methods are much more robust between the calibration and

$RB_{MB}$  versus  $RB_O$  for the temperature indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.



**Figure 3.** Box plot of the Annual number of days with precipitation higher than 10 mm (ETCCDI 'Heavy precipitation' days, see Zhang et al. (2011)) in the validation period. (a) observations, (b) raw simulations, (c) QDM, (d) mQDM, (e) MBCn, (f) MRQNBC, (g) dOTC, (h)  $R^2D^2$ .

700 validation period: the performance of the univariate methods is worse in all seasons. Nonetheless, the univariate methods still perform better.

#### 4.4 Potential evaporation

**Table 4.** PSS values for temperature in the calibration (Cal) and validation (Val) periods (%).

|          | Winter |      | Spring |      | Summer |       | Autumn |      |
|----------|--------|------|--------|------|--------|-------|--------|------|
|          | Cal    | Val  | Cal    | Val  | Cal    | Val   | Cal    | Val  |
| QDM      | 97.1   | 93.4 | 95.8   | 86.8 | 96.8   | 88.7  | 96.5   | 91.4 |
| mQDM     | 99.3   | 94.0 | 98.8   | 87.0 | 99.1   | 89.8  | 98.7   | 91.8 |
| MBCn     | 52.7   | 52.5 | 78.6   | 77.1 | 44.8   | 39.2  | 77.9   | 79.3 |
| MRQNBC   | 78.7   | 76.6 | 90.6   | 75.3 | 58.9   | 61.7  | 87.1   | 80.7 |
| dOTC     | 81.6   | 82.0 | 81.8   | 83.0 | 79.2   | 77.5  | 80.3   | 83.3 |
| $R^2D^2$ | 71.7   | 69.7 | 75.5   | 72.2 | 63.6   | 59.22 | 73.4   | 73.2 |

Figure ?? displays the Although the PDF of the adjusted simulations matches the observed PDF relatively well, the  $RB_{MB}$  and  $RB_O$  values (Figure 4) show some clear differences between the seasonal bias adjustment: for winter (blue) all methods

705 perform poorly, whereas for the other seasons, at least some methods are able to adjust the raw simulations. For winter, the  
R index values are high for all percentiles, which indicates that nonstationarity could be the cause for the poor performance.  
However, this is not clear-cut. When comparing the winter  $RB_{OMB}$  and  $RB_{MB}$  values for the E indices. Only a few indices  
are shown for each method, or just one for dOTC of the validation period with those of the calibration period (Fig. 5; blue).  
710 it can be seen that only QDM (panel (a)) performs much better and that mQDM, MRQNBC and dOTC (respectively panels  
(b), (d) and (e)) perform slightly better in the calibration period. The better performance of these methods is clearest for the  
lower percentiles ( $T_5$ ,  $T_{25}$  and  $T_{50}$ ). MBCn (panel (c)) and  $R^2D^2$  (panel (f)) seem to perform equally poor in both calibration  
and validation period. The poor performance of these two methods could be caused by the seasonal evaluation: both apply a  
shuffling algorithm over the full time period. However, indicating that the performance after adjusting the bias is generally  
worse than the raw climate simulations. The indices plotted are  $E_{25}$ ,  $E_{99}$  and  $E_{99.5}$ . The index  $E_5$  also performs well, but cannot  
715 be plotted as its observed value is 0 mm. Thus, for the lowest and the highest percentiles, the bias-adjusting methods perform  
well, but they fail to capture the nonstationarity at the middle percentiles. These middle percentiles have high R index values:  
they are all greater than or equal to one. Only for dOTC, it is possible to plot a percentile between  $E_{25}$  and  $E_{99}$ :  $E_{95}$  other  
methods, this is harder to explain: QDM, mQDM and MRQNBC all use seasonal time windows, while dOTC does not. How-  
ever, for dOTC, this QDM and mQDM, the moving time window used in the adjustment and the fixed seasonal window in  
720 the evaluation might cause a small mismatch. For MRQNBC, there is also the only percentile for which it is possible to  
plot the  $RB_O$  and  $RB_{MB}$  influence of the monthly and yearly adjustment. For dOTC, the optimal transport and, hence, stochastic  
element might be better suited for seasonal differences than the shuffling used by MBCn and  $R^2D^2$ , but still does not seem  
optimal. Besides, the seasonal variance is larger for temperature than for precipitation, which increases the susceptibility of the  
methods to differences in seasonal adjustment and evaluation. As a last reason, it should be considered that the  $RB_{MB}$  values  
725 (respectively 1.00 and 0.86). For all other methods and all other percentiles, both  $RB_O$  and  $RB_{MB}$  values are higher than 1.  
The poor performance of dOTC might be related to its trend-preservation characteristics. Of all methods, it is the one most  
explicitly designed to follow the simulated trend. This might thus imply that the nonstationarity for E is caused by poor model  
performance, although this should be investigated more in depth. and  $RB_O$  values always depend on respectively the original  
biases and the observations.

730 The percentiles that are plotted all have a high  $RB_O$  value, which is in this case caused by rather low biases to adjust. For  
example,  $E_{99.5}$  had an observed value of 5.24 mm/day and only a bias of 0.27 mm/day, or 5%. There is consequently not much  
room for improvement, though the  $RB_{MB}$  values imply that the bias-adjusting method could be improved, except for  $E_{99.5}$ ,  
which consistently has an  $RB_{MB}$  value lower than 0.5.

$RB_{MB}$  versus  $RB_O$  for the potential evaporation indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

735 As in Section 4.3, there are interesting differences between QDM and mQDM. In contrast to the results for temperature,  
mQDM performs better. For mQDM, At first sight, in spring (ochre), most methods, with the exception of MBCn (panel (c))  
and MRQNBC (panel (d)), seem to perform relatively well. However, when comparing the biases of the validation period with  
those of the calibration period, the adjustment of  $T_5$  by QDM, mQDM, MRQNBC and dOTC (respectively panels (a), (b), (d)  
and (e)) is clearly poorer, whereas the highest percentiles ( $ET_{99}$  and  $ET_{99.5}$ ) have lower  $RB_{MB}$  values than for QDM. In this

740 ease, it thus seems better to only use the climate change signal. However, given the general poor performance of all methods, these results should be considered with care.

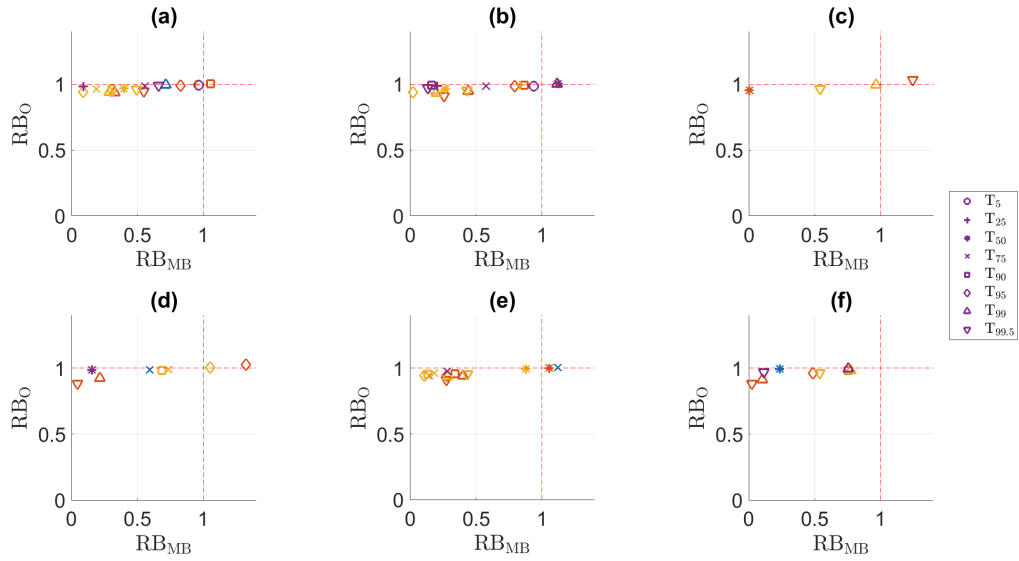
Given that the percentiles with a perform similar to the calibration period or better. For MBCn and  $R^2D^2$ , the performance is similarly. The poor performance corresponds to the high R index value have a larger bias than the raw simulations after adjustment, and the added value for the other percentiles with respect to observed values is low, it can be advised not to adjust E. However, similar to for T, this should be evaluated on a case-by-case basis.

#### 4.4 Correlation

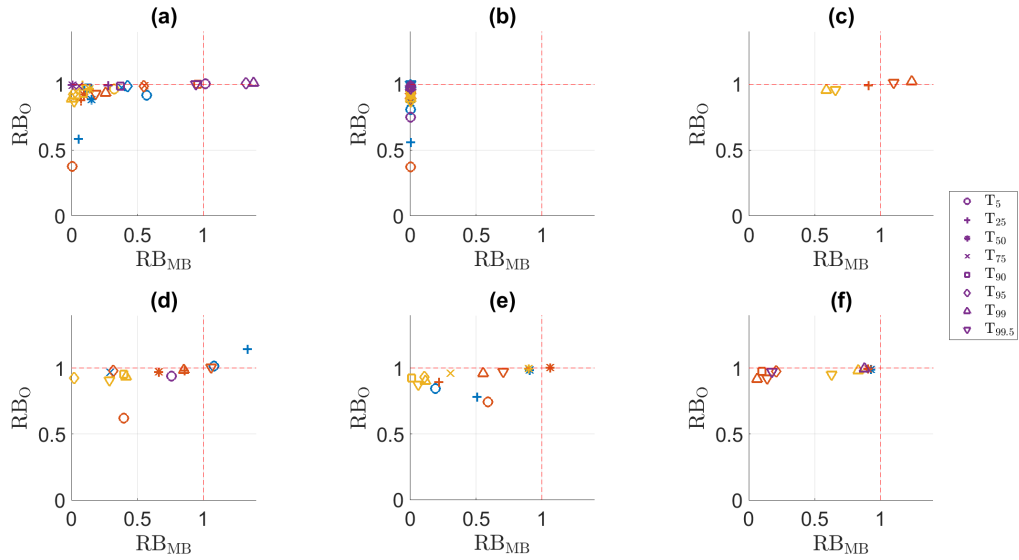
When considering the correlation (Fig. ??), the methods generally perform well: most of the correlation indices can be plotted. For summer, this is also observed, although to a smaller extent: only QDM and mQDM were able to properly adjust  $T_5$  in the calibration period. In general, QDM, MRQNBC and dOTC all perform slightly worse in the validation period in comparison with the calibration period for summer. mQDM performs similarly, whereas MBCn and  $R^2D^2$  perform poorly in both periods. In autumn, the performance is poor for all methods in the validation period. However, there are some differences depending on the indices under consideration and the method. The indices that can always be plotted are the lag-0 cross-correlation between P and T, the lag-1 cross-correlation between P and T, the lag-1 cross-correlation between P and E and the correlation between P and T. Except for dOTC, all methods also perform well for the lag-0 cross-correlation between P and E. Yet, for the indices that can be plotted, the  $RB_O$  and  $RB_{MB}$  values show a considerable difference among the methods. For example, except for QDM, the performance is poor in the calibration period as well, and, hence, conclusions are hard to draw. However, based on the R index values, which indicate limited nonstationarity, it could be assumed that the influence of the seasonality is larger than that of the nonstationarity.

Based on the results for winter and the lag-1 cross-correlation between P and E has  $RB_O$  values ranging from 0.29 (mQDM) to 0.69 (dOTC) and  $RB_{MB}$  values ranging from 0.08 (mQDM) to 0.60 (dOTC). The best method varies for each index: while dOTC does not perform well for most indices, it has the best performance for the correlation between P and T. It is interesting that the performance of dOTC seems to be either very good, or very poor. As dOTC is built around the idea of trend preservation and all-in-one adjustment, it is possible that the adjustment performs very well when the trend is properly modelled. Three out of the four correlations that dOTC adjusts well are based on T and P, the two variables that are well understood in the time frame under consideration. All indices that are not or less frequently present in the plots have one thing in common: they are based on the correlation between E and one of the other variables. The indices based on E thus consequently perform worst, except for the aforementioned lag-1 cross-correlation between P and E. lowest percentiles in spring and summer, it seems that the lower temperature values are more susceptible nonstationarity. This should certainly be accounted for when estimating extremes such as cold spells.

770 The correlation index performance seems to be related to the results of T (Section 4.3) and E (Section 4.4): correlations of E with another variable generally perform worse, and the (cross-)correlation between T and E with another variable performs the worst, in line with the R index values (Table ??). Although the



**Figure 4.**  $RB_{MB}$  versus  $RB_0$  for the correlation-temperature indices in the validation period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple..



**Figure 5.**  $RB_{MB}$  versus  $RB_0$  for the temperature indices in the calibration period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.

#### 4.4 Potential evaporation

775 The PSS values for potential evaporation (Table 5) show that the univariate bias-adjusting methods perform better than the mul-  
 780 tivariate bias-adjusting methods. They are supposed to adjust correlation, they seem to be unable to do so, as they generally have larger  
 biases (though not on all indices) than the univariate bias-adjusting methods for the indices with the lowest residual bias values.  
 This seems to indicate that the multivariate bias-adjusting methods, and especially MBCn and dOTC, are unable to adjust the  
 correlation exactly because of the nonstationarity in the correlation that has to be overcome. In contrast when considering the  
 full PDF. Similarly to temperature (Table 4) the skill scores differ among the multivariate methods. However, in contrast to  
 785 temperature, dOTC performs much worse for potential evaporation; MRQNBC performs best. Similarly to temperature, MBCn  
 and  $R^2D^2$  vary depending on the season. In comparison with the calibration period, the univariate bias-adjusting methods  
 neglect the adjustment of correlation and consequently do not have to overcome nonstationarity in the correlation bias. Yet,  
 for the univariate bias-adjusting methods, the difference in adjustment of T and E seems to have an influence here as well, as  
 illustrated by the different results for QDM and mQDM. Except for the correlations that have high  $RB_{MB}$  values for all meth-  
 790 ods perform worse in the validation period in every season, whereas this varies for the multivariate methods: only in spring  
 and summer, all multivariate methods perform worse. For spring, the difference is large, which could be related to the clear  
 nonstationarity for this season. For summer, the R index values are generally lower, which indicates less nonstationarity, but  
 the difference in PSS between calibration and validation period is also smaller. The large difference for spring between both  
 periods is striking, as this was not as apparent for winter temperatures, despite the high R index values. This could be explained  
 795 by the R index values for the results indicate that mQDM performs better. Thus, it might be better to use correlations of the  
 observed time series than to adjust the simulated correlations. This is confirmed by the results for MRQNBC. Together with  
 mQDM, this is the only method to have six indices with  $RB_{MB}$  values lower than one. Besides, it is also the only multivariate  
 bias-adjusting method to have an  $RB_{MB}$  value for  $corr_{E,T}$  lower than one, although it is only slightly lower. mean and standard  
 deviation: for potential evaporation in spring, only the bias in the mean changed a lot, whereas for temperature in winter, both  
 the biases in mean and standard deviation changed a lot. The combination of these bias changes could offset each other in the  
 calculation of the PSS.

#### 4.5 Precipitation occurrence

Table 5. PSS values for evaporation in the calibration (Cal) and validation (Val) periods (%).

|          | Winter |      | Spring |      | Summer |      | Autumn |      |
|----------|--------|------|--------|------|--------|------|--------|------|
|          | Cal    | Val  | Cal    | Val  | Cal    | Val  | Cal    | Val  |
| QDM      | 94.8   | 86.1 | 93.3   | 82.5 | 97.3   | 88.6 | 94.1   | 91.0 |
| mQDM     | 100.0  | 90.5 | 100.0  | 83.4 | 100.0  | 87.7 | 100.0  | 92.4 |
| MBCn     | 48.5   | 52.0 | 78.6   | 70.8 | 52.7   | 48.4 | 79.5   | 83.6 |
| MRQNBC   | 89.1   | 84.5 | 91.4   | 74.1 | 80.2   | 78.0 | 85.1   | 88.6 |
| dOTC     | 58.7   | 52.4 | 67.4   | 57.5 | 63.9   | 56.0 | 60.5   | 57.0 |
| $R^2D^2$ | 80.4   | 79.3 | 69.1   | 59.5 | 66.5   | 63.9 | 78.6   | 76.5 |

Figure ?? displays the  $RB_{MB}$  and  $RB_O$  and  $RB_{MB}$  results for potential evaporation in the validation period are displayed in Figure 6. For every season, all methods perform rather poorly, although there are differences between the method's performances and in the extent of nonstationarity. Based on the R index values and Table 5, it would seem that spring is most influenced by bias nonstationarity, as many percentiles have an R index value higher than 1 and the PSS values differ considerably for spring. Figure 6 shows that only  $E_5$  (for QDM, mQDM, MRQNBC and  $R^2D^2$ , respectively panels (a), (b), (d) and (f)),  $E_{99}$  (for mQDM and MBCn, panels (b) and (c)) and  $E_{99.5}$  (for mQDM and MBCn, panels (b) and (c)) have  $RB_{MB}$  and  $RB_O$  values for lower than 1. Except for  $E_{99.5}$ , this corresponds to the precipitation occurrence indices. When comparing these values, large differences among the methods and among the indices can be noted. The best-performing method seems to be QDM, as all the  $RB_{MB}$  values are lower than 0.5 and some  $RB_O$  values are close to 0.5. When comparing the other methods percentiles that have an R index value lower than 1. For mQDM and MBCn, there is no clear difference between the univariate and multivariate bias-adjusting methods. The other univariate method, mQDM, and one multivariate method, MRQNBC, also perform better than it cannot be ruled out that the good performance for  $E_{99.5}$  is by accident. However, bias nonstationarity alone does not explain the poor performance: when comparing the biases in the calibration (Fig. 7) and validation periods, it can be seen that, except for mQDM, all methods perform poorly in the calibration period. For MBCn, dOTC and  $R^2D^2$ , which perform the worst in the raw climate simulations for all indices. The other two methods, MBCn and dOTC, have respectively only one and two indices with both  $RB_O$  and  $RB_{MB}$  values below 1. calibration period, this could be related to the absence of a seasonal component, whereas this is less clear for QDM, mQDM and MRQNBC, as discussed in Section 4.3. Nonetheless, the latter three methods are all able to adjust  $E_{25}$  and  $E_{50}$ , two percentiles that cannot be adjusted by any method in the validation period.

$RB_{MB}$  versus  $RB_O$  for the precipitation occurrence indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC. Interestingly enough, the indices with  $RB_O$  and  $RB_{MB}$  values below 1 are not the same for all methods. For the three best-performing methods, in the other seasons, the dry-to-dry transition probability has very low  $RB_{MB}$  values (ranging from 0.3 methods behave similarly to spring: most of the multivariate methods perform as poorly in the calibration period as in the validation period (except MRQNBC, to 0.18), while this index is absent from the MBCn and dOTC plots. The differences between those two plots are also notable. For MBCn, only the number of dry days has a very low  $RB_{MB}$  value (0), as the number of dry days is unaffected after the thresholding, whereas the lag-1 P auto-correlation and some extent). The poor performance of the multivariate methods in the calibration period indicates that the absence of a seasonal component might have a large impact, as was also discussed in Section 4.3. This is confirmed by the results for the full year (not shown), which show that all methods perform well in the calibration period.

Despite the poor performance of some methods in the calibration period, even for these seasons some differences between the calibration and validation period are worth discussing. In winter (blue), where nonstationarity mostly affected the standard deviation, the performance of all methods for all indices is slightly worse in comparison with the calibration period. Only the lower percentiles ( $E_5$  and  $E_{25}$ ) can be adjusted well by almost every method. In summer (yellow), where the R index values indicated some nonstationarity for the lower E percentiles, the performance is poorer in the validation period for all percentiles except  $E_{99}$  and  $E_{99.5}$  (and  $E_{90}$  for dOTC). However, the impact seems to be smaller for MBCn, dOTC and  $R^2D^2$ . In autumn

(purple), the R index values indicated the largest impact on the standard deviation. As in winter, the wet-to-dry transition probability are more biased than the raw climate simulations. For dOTC, it is the other way around: the number of dry days is more biased after the application of dOTC, and the auto-correlation and the wet-to-dry transition probability perform well, with  $RB_{MB}$  values 0.50 or lower. best performance is obtained for the lowest percentiles and, for the univariate methods, for the highest percentiles ( $E_{99}$  and  $E_{99.5}$ ). Despite the seemingly larger impact on the univariate methods in these three seasons, their adjustment is still better than the adjustment by the multivariate methods.

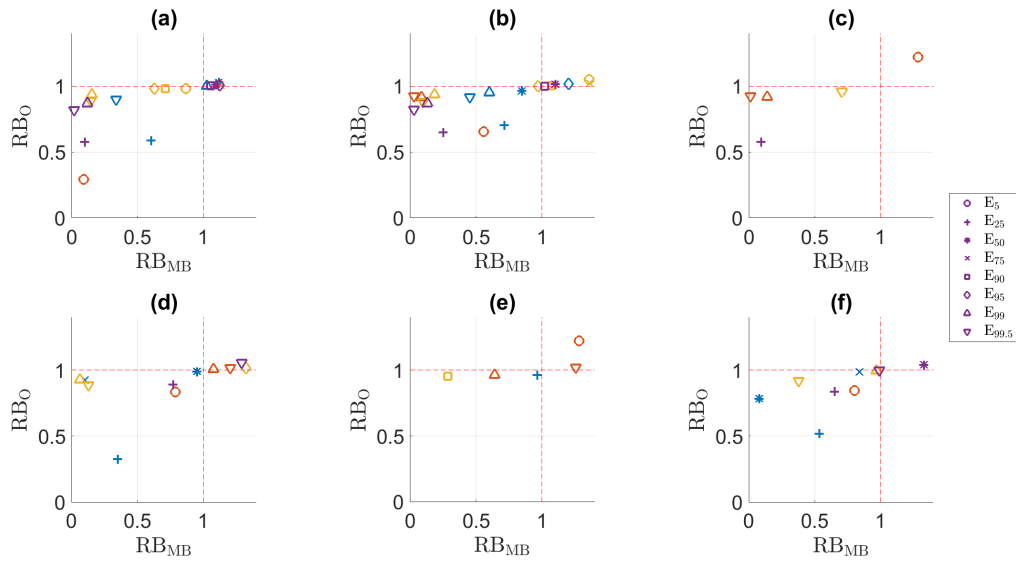
Another peculiar result that can be seen from Fig. ?? is the difference in dry day bias. Although all methods start with the same number of dry days, there are large differences among the  $RB_O$  and  $RB_{MB}$  values. The results for potential evaporation have to be considered in comparison to the effective bias values for the original simulations and the adjusted simulations: the original biases were relatively small (not shown). Hence, even a small change in bias will have a large impact. Nonetheless, even these small changes and relatively small biases have an impact, which is reflected by the  $RB_{MB}$  values for the number of dry days. The  $RB_O$  values range from 0.58 to 1.04 and the  $RB_{MB}$  values from 0 to 1.09. QDM and MBCn perform best ( $RB_{MB} = 0$ ), as the number of dry days is unaffected after the thresholding. For mQDM ( $RB_{MB} = 0.25$ ), this holds by construction: instead of adjusting the threshold-adjusted climate model simulations, this method adjusts the observations. For MRQNBC ( $RB_{MB} = 0.63$ ) and dOTC ( $RB_{MB} = 1.09$ ), the results seem to imply that the multivariate framework of these methods has an influence on the number of dry days values. On the other hand, when considering the PSS values, which reflect the full PDF instead of focusing on the extremes, the impact is limited, although this depends on the method and season, as was shown for spring.

What the difference in transition probabilities implies for the time series, becomes more clear in Fig. ?. Although all adjusted simulations and the observations have more short wet spells than long ones, MBCn pronounces the short wet spell length more than the other methods, while the probability of longer wet spell lengths is lowered in comparison with other methods. Closest to the observations is mQDM. QDM and MRQNBC also perform well, a conclusion similar to that of Fig. ?.

The difference in performance between QDM and the other methods seems to demonstrate that most strategies for retaining a certain temporal structure or adjusting the temporal structure do not perform well. MRQNBC and mQDM depend heavily on the temporal structure of the observations, and MBCn and dOTC have an important shuffling or recalculation aspect, all of which lead to less reliable results at the end of the process. The poor performance of dOTC and MBCn for the temporal structure was also discussed by François et al. (2020). As for mQDM and MRQNBC, it is notable that the temporal structure does not change much from the calibration time series to the validation time series. At least, this is suggested by their relatively good performance, which is based on using the observed time series (mQDM) and observed persistence statistics (MRQNBC). Yet, this is no guarantee that these methods will be able to realistically adjust climate model simulations for the end of the century.

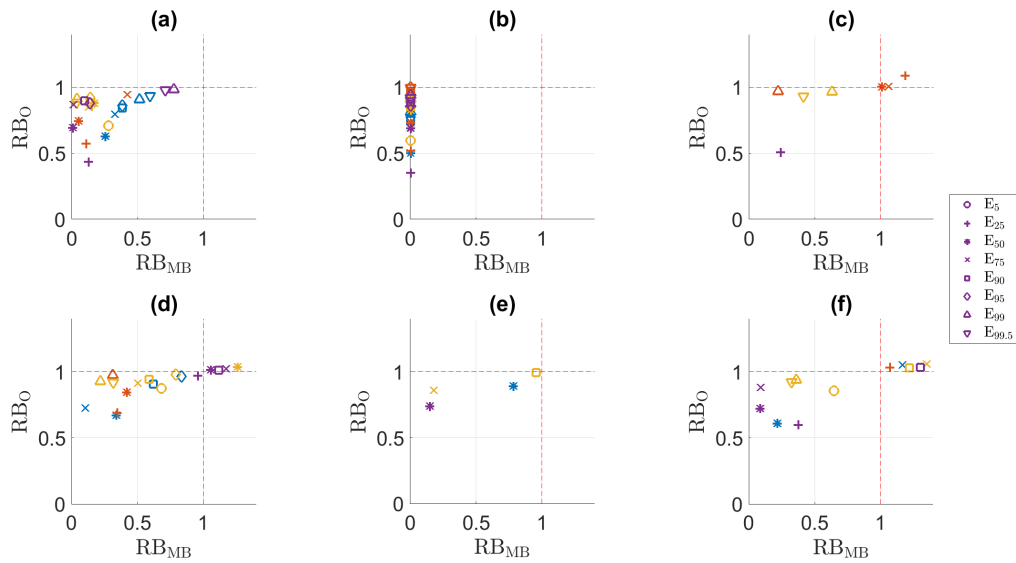
Figure ?? also suggests that despite the high R index value, the P-lag-1 auto-correlation is not necessarily poorly adjusted. For QDM, this index has relatively low  $RB_O$  and  $RB_{MB}$  values. This could imply that the performance still depends on the robustness to the bias nonstationarity of the methods under consideration. Or, as the other indices illustrate, the effect of bias (non-)stationarity is not as large as the effect induced by the methods themselves. An example of this is the number of dry days: though it has a low R index value, the performance varies substantially among the methods.





**Figure 6.**  $RB_{MB}$  versus  $RB_0$  for the potential evaporation indices in the validation period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.

Wet spell length probability mass function for all adjusted simulations, the raw RCM simulations and the observations.



**Figure 7.**  $RB_{MB}$  versus  $RB_0$  for the potential evaporation indices in the calibration period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f) Winter: blue, spring: ochre, summer: yellow, autumn: purple.

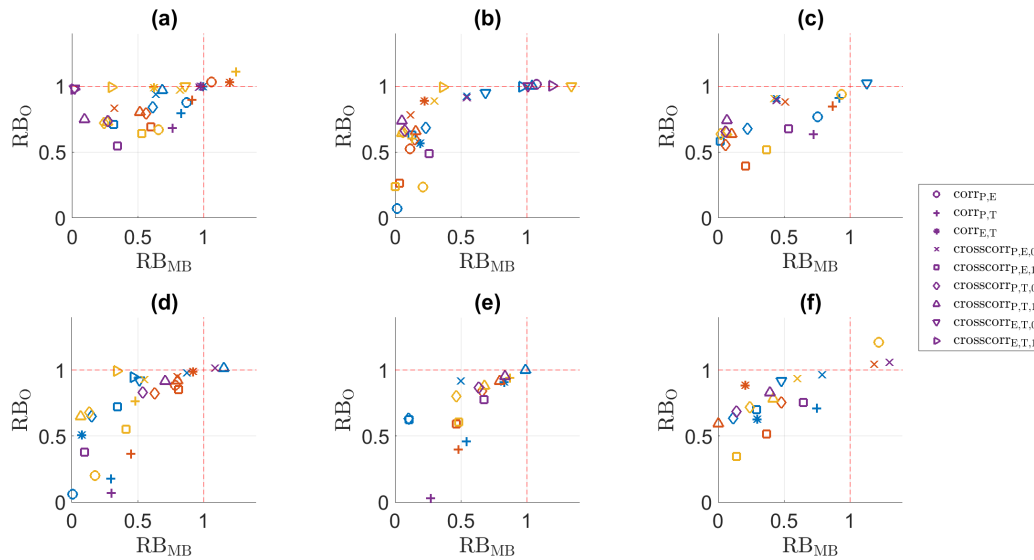
## 4.5 Correlation

## 4.6 **Discharge**

870 All bias-adjusting methods perform better for the discharge percentiles compared to most other indices (Fig. ??). Although the discharge is influenced by a combination of many effects, these appear to be small in the end result. For example, the poor performance of E-For correlation (Fig. ??) does not result in large discharge biases. Thus, it is the integration of precipitation amount, precipitation occurrence and evaporation, and the routing effect that ultimately defines the resulting discharge. Generally, many indices perform well, though there are some differences among the methods. The best-performing methods are QDM and mQDM, as all indices have values lower than 0.5 and some indices have values near to zero. For mQDM, 875 8), all methods perform relatively well in the validation period. Both the univariate and the 20-year return period even has an  $RB_O$  value of -0.05. For MBCn, the results are also good. Most extreme values have  $RB_O$  and  $RB_{MB}$  values lower than 0.5; only the 5th percentile has an  $RB_O$  value of 0.96 and an  $RB_{MB}$  value of 0.89. As the only difference between QDM and MBCn is the adjustment of occurrence, the results for discharge illustrates the importance of occurrence adjustment. This variability in values seems to be a difference between the univariate-multivariate bias-adjusting methods can adjust the simulated biases well. The univariate methods will adopt the dependence structure of the raw simulations, whereas the multivariate methods are specifically designed to adjust the dependence structure, and both strategies seem to work well. However, it should be noted that some of the biases in correlation are very small in the raw simulations (not shown) and that for those correlations, the good adjustment by univariate methods is trivial: they will adopt the correlation of the simulations and only slightly adjust this by adjusting the marginals. This is linked with an issue raised by Zscheischler et al. (2019): in situations with low biases in the correlation, the univariate methods will almost always outperform the multivariate bias-adjusting methods. For the two worst performing methods, i. e. MRQNBC and dOTC, some indices have  $RB_O$  and  $RB_{MB}$  values close to or, as specifically adjusting the dependence structure sometimes results in an increase of the bias.

890 The good performance for the validation period indicates that the impact of nonstationarity is limited, as was also shown by the small R index values (Section 4.1). This is confirmed by the biases in the calibration period (not shown), which are similar to those in the validation period. However, for some values, the R index value was higher than 1 and some values between 0.5, thus it is important to know what caused this. For  $corr_{E,T}$  in summer, the difference between the validation and calibration period is negligible, although only for QDM this value is well adjusted in both periods. However, the bias for the original simulations is lower than 0.10% in both the calibration and validation period, and 1. These methods are thus unable to correctly adjust the bias for all indices. However, although MRQNCB and dOTC seem to perform similarly, the indices with the worst  $RB_O$  and  $RB_{MB}$  values are different. For MRQNBC, the 99.5th percentile and the 20-year return period have the highest values, whereas for dOTC, the 5th and 25th percentile perform worse than the raw climate simulations. From the point of view of extreme discharges, dOTC is thus the better method of these two. This might indicate that although not all occurrence indices of dOTC had lower  $RB_{MB}$  values than those of MRQNBC, those that had ( $P_{lag1}$  and  $P_{PT0}$ ), had a larger influence on the extreme discharge values. Both indices are partly based on the occurrence of wet days, and thus indicate that those need to be at the correct place in the time series for extreme floods to be correctly simulated switches in sign, which

inflates the R index value. For  $\text{crosscorr}_{E,T,0}$  and  $\text{crosscorr}_{E,T,1}$ , the same effect occurs. Besides, it seems that the bias of these three correlations is too small to be corrected by any method and that trying to adjust this automatically inflates the results. As discussed earlier, this shows that while the R index can be a valuable tool for some variables, it does not always tell the full story.



**Figure 8.**  $RB_{MB}$  versus  $RB_0$  for the discharge percentiles and correlation indices in the 20-year return validation period value. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.

## 5 Discussion

### 4.1 Precipitation occurrence

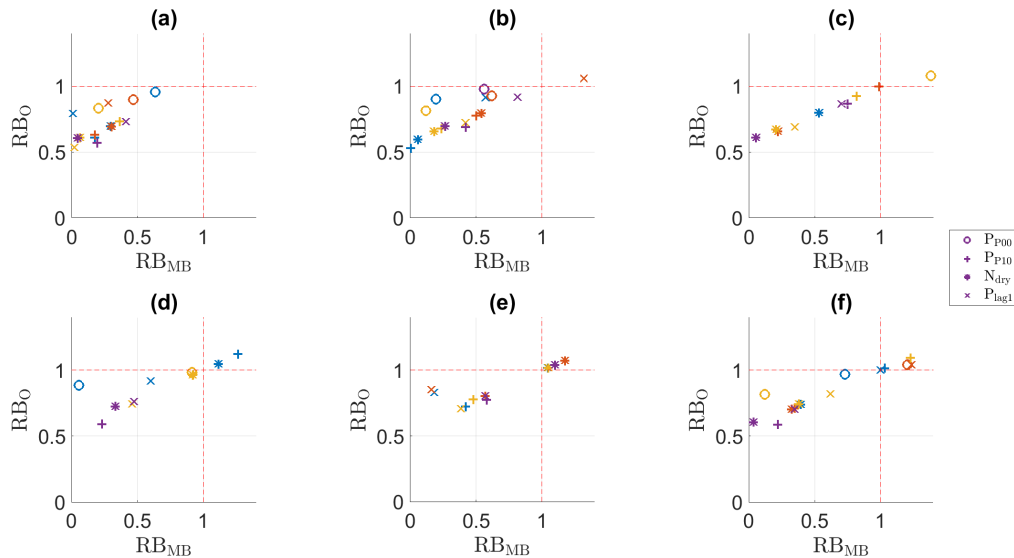
In the previous section, the results for the bias adjustment by different methods and under climate change conditions were reported. The effect of climate change on the bias was evaluated through the R index, which showed that the bias for some indices cannot be considered stationary. For some of the indices (the lower percentiles of T and especially the middle percentiles of E) the methods performed poorly, which could often be linked with the Figure 9 shows that the bias-adjusting methods are able to adjust the precipitation occurrence well in most seasons. Especially the univariate bias-adjusting methods perform well. Although the multivariate bias-adjustment always results in at least one index that is better than the raw climate simulations (except for MRQNBC in spring: panel (d), ochre), most indices are not, or only slightly better than the raw climate simulations. This is a disadvantage inherent to the current generation of multivariate bias-adjusting methods: as discussed in Section 3.3, the dependence adjustment will always influence the temporal structure (François et al., 2020; Vrac and Thao, 2020b). Nonetheless,

on a seasonal level, the temporal structure is sometimes remarkably well adjusted, such as in summer (yellow) and autumn (purple).

920 The R index values. The methods clearly handle this bias nonstationarity differently. It seems that the univariate bias-adjusting methods are far more robust: even for indices with high R values, they are sometimes able to perform very well, with low  $RB_{\text{O}}$  and  $RB_{\text{MB}}$  values. This good performance thus seems to imply that the more indices a bias-adjusting method directly adjusts, the more susceptible it is to problems related to bias nonstationarity. However, indicated that there might be some nonstationarity in spring and autumn (Section 4.1): the value for  $P_{\text{lag}1}$  is 2, and for the other indices the values are clearly higher than those in winter and summer. In contrast to other situations of bias nonstationarity, this does not imply that QDM and mQDM are  
925 similar: while they are almost as good for many variables, the poorer performance of mQDM for the precipitation occurrence indices is an indication that assuming that the temporal structure of the past can be used for the future might be dangerous, as Johnson and Sharma (2011) and Kerkhoff et al. (2014) already mentioned. Given that mQDM performed worse for two time periods separated by 10 years only, it is unlikely that it is safe to use this method, or other delta change-based methods, for impact assessments targeting the end of the 21st century and depending on the temporal structure of time series. Yet, for some  
930 other indices, especially the correlation, mQDM performed better. Consequently, the exact choice should depend on the goals of the end user. result in a poorer, but actually a better performance for these two seasons (calibration period not shown). Winter and summer, for which no nonstationarity could be detected, perform similarly in both the calibration and validation period. However, in all seasons mQDM (panel (b)) performs worse in the validation than in the calibration period. As this method uses the observed structure, the temporal structure is by construction perfect in the calibration period. The poorer result in the  
935 validation period might imply that using the observed temporal structure does not suffice for future impacts, which might be important when using delta methods for impact assessment.

The results of the multivariate bias-adjusting methods too are not without nuance: though they are generally worse than the univariate bias-adjusting methods, their performance depends heavily on the variable under consideration and on the method itself. A clear example of this dependence on variables is the contrasting performance of dOTC to adjust T (Fig. ??) versus  
940 E (When comparing the methods, some differences related to their structure can be noticed. In general, QDM (panel (a) in Fig. ??): the adjustment of E is much worse. This is a reminder that in a multivariate context, the multivariate methods are far less robust and can perform relatively good and poor at the same time for different variables. Therefore, there seems to be an interplay between the modelling of the variables and the method of calculation. Except for P, for which the results were similar, the methods performed differently for each variable. MRQNBC performed best in the context of temporal structure, for  
945 which it was designed (Fig. ??). For T, MBCn and dOTC performed better (Fig. ??) 9) has the best performance of all methods for the occurrence, indicating once more the impact of the shuffling and similar algorithms of the multivariate bias-adjusting methods. This could be related to their trend-preserving properties, which are more pronounced for those methods than for MRQNBC. For E (Fig. ??) and correlation (Fig. ??) Only in autumn (purple), dOTC displayed the most different results. For the former, MRQNBC (panel (d)) and  $R^2D^2$  (panel (f)) perform as well as QDM. However, mQDM (panel (b)) also performs  
950 well in all seasons, despite the poorer fit. There are also differences among the different multivariate bias-adjusting methods. In all seasons, MBCn (panel (c)) and  $R^2D^2$  (panel (f)) are able to reduce the bias of the number of dry days, whereas this

955 varies for MRQNBC and dOTC (panel (e)). The good performance for this index for MBCn and R<sup>2</sup>D<sup>2</sup> is based on the use of thresholding and QDM for the marginal adjustment: these methods are able to perfectly adjust the number of dry days, and any remaining bias can be related to the combination of temporal shuffling and seasonal evaluation. However, dOTC adjusts P<sub>lag1</sub> and P<sub>P10</sub> well in every season. This implies that it is able to differentiate in the adjustment between zero and non-zero values, whereas longer series of zeros are harder to adjust. The incorrect series of zeros is probably also linked with one of the deficiencies of dOTC: it sometimes creates nonphysical precipitation values, which have to be corrected by thresholding.



**Figure 9.**  $RB_{MB}$  versus  $RB_0$  for the precipitation occurrence indices in the validation period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f) R<sup>2</sup>D<sup>2</sup>. Winter: blue, spring: ochre, summer: yellow, autumn: purple.

## 4.2 Discharge

960 The Perkins Skill Score values for discharge (Table 6) show that the univariate bias-adjusting methods generally perform best, whereas the performance of the all-in-one and trend-preservation method did not seem robust enough. For the latter, it depended heavily on the type of correlation under consideration. These results seem to imply that the difference under bias-nonstationary conditions is not clear-cut for the different types of multivariate bias-adjusting methods. For the ‘marginal/dependence’ vs. the ‘all-in-one’ approach, consequently no clear conclusions can be drawn. For the amount of temporal alteration, it depends on the index under consideration. MRQNBC, which replaces the simulated correlations by those of the observations performs well for the temporal structure, but performs worse for many other indices. For MBCn and dOTC, the effect of the difference in temporal alteration is less distinct and other properties, such as trend preservation, seem to have more influence depends on the season. However, all methods perform poorly for spring. The PSS values for evaporation clearly show the impact of nonstationarity, which seems to be propagating to the discharge PDF. This is illustrated when comparing with the PSS values

965

970 for the calibration period: only in spring, all methods perform worse in the validation period than in the calibration period. For the other seasons, the impact is much more mixed.

To have a better view of how these results should be interpreted, the perspective of the end user should be considered (Maraun et al., 2015; Maraun and Widmann, 2018b). We used discharge as an example, using the relatively simple PDM. Although the residual bias values

**Table 6.** PSS values for discharge in the calibration (Cal) and validation (Val) periods (%).

|                               | Winter |      | Spring |      | Summer |      | Autumn |      |
|-------------------------------|--------|------|--------|------|--------|------|--------|------|
|                               | Cal    | Val  | Cal    | Val  | Cal    | Val  | Cal    | Val  |
| QDM                           | 85.9   | 90.0 | 87.4   | 67.0 | 90.8   | 81.5 | 92.4   | 86.6 |
| mQDM                          | 99.6   | 85.3 | 100.0  | 60.8 | 100.0  | 76.8 | 100.0  | 86.2 |
| MBCn                          | 49.5   | 50.7 | 87.4   | 51.0 | 41.8   | 64.3 | 74.3   | 80.0 |
| MRONBC                        | 92.2   | 86.4 | 67.3   | 38.6 | 89.1   | 85.1 | 57.7   | 49.8 |
| dOTC                          | 70.0   | 85.4 | 48.2   | 25.0 | 42.3   | 58.3 | 68.5   | 72.4 |
| R <sup>2</sup> D <sup>2</sup> | 75.0   | 78.0 | 73.1   | 42.8 | 43.2   | 40.0 | 69.3   | 63.8 |

975 The impact on the PDF for spring discharge does not clearly appear when comparing the  $RB_{MB}$  and  $RB_0$  values: for all methods for E (Fig. ??) indicate a poor performance, the influence thereof on the discharge seems to be negligible and seasons, the bias adjustment seems to result in an agreeable representation of the discharge in the validation period (Fig. ??). Discharge is the variable that is of the highest importance for hydrological impact modelling, and the results indicate that most methods are able to adjust the forcing variables sufficiently in order to have a good simulation of discharge. However, the small differences between the methods should still be taken into account. Overall, QDM and mQDM perform best in adjusting the variables such that the discharge rates are the least biased in comparison with the observations. This is also important considering that bias adjustment can be applied for many different types of impact assessment. In other impact assessments, (10). However, when comparing these results with the residual biases in the differences could affect the result more than the discharge considered here. For example, forest fires (a typical compound event, discussed in a bias adjustment context in e.g. Yang et al. (2015); Cannon (2018); Zscheischler et al. (2019)) depend more heavily on T and E to simulate fire weather conditions. Besides such compound events, other applications are ecosystem functioning (Sippel et al., 2016), agriculture (Galmarini et al., 2019), or climate zone classification (Beck et al., 2018). In such studies the effect of bias nonstationarity can even be worse, whereas in other studies, depending more on P or other (so far) less-affected variables, the need for a bias nonstationarity-proof calibration period (Fig. 11), it becomes clear that the results for winter and summer are much worse in the validation period. This corresponds with the poor performance for precipitation adjustment in these seasons, which was probably linked with bias nonstationarity.

985 The bias-adjusting method is less compelling. Anyway, methods seem to respond similarly to the nonstationarity. In winter (blue), QDM (panel (a)) performs slightly better, whereas in summer (yellow), R<sup>2</sup>D<sup>2</sup> (panel (f)) performs relatively good. In spring (ochre), the inability of some methods to adjust the biases in nonstationary conditions implies that a thorough

assessment of possible bias nonstationarity should be made before bias adjusting. If not done, the risk of reporting a wrong  
995 future projection is likely increased. Given the knowledge of bias nonstationarity, such uncertainties can be better characterised.  
methods also perform similarly, although QDM performs slightly better for  $Q_{99}$  and dOTC (panel (e)) performs worse than the  
other methods. As such, whether or not the methods take seasonality explicitly into account does not seem to matter for the  
impact on discharge. This also follows from the structure of the hydrological model: precipitation is a more important driver  
than potential evaporation. Seasonality in the bias adjustment had a larger impact on potential evaporation, but this impact  
1000 disappears when using these variables as inputs to the hydrological model. Besides, it can also be seen that if an important  
forcing variable for an impact model shows large nonstationarity, this nonstationarity will propagate through the model. This  
helps explaining the differences between the PSS and the RB values: the impact of nonstationarity on potential evaporation  
propagates as an influence on the PDF structure, but is less visible in the final bias, as the amount of precipitation has a much  
larger impact in the hydrological model. Hence, the final bias is more influenced by precipitation nonstationarity.

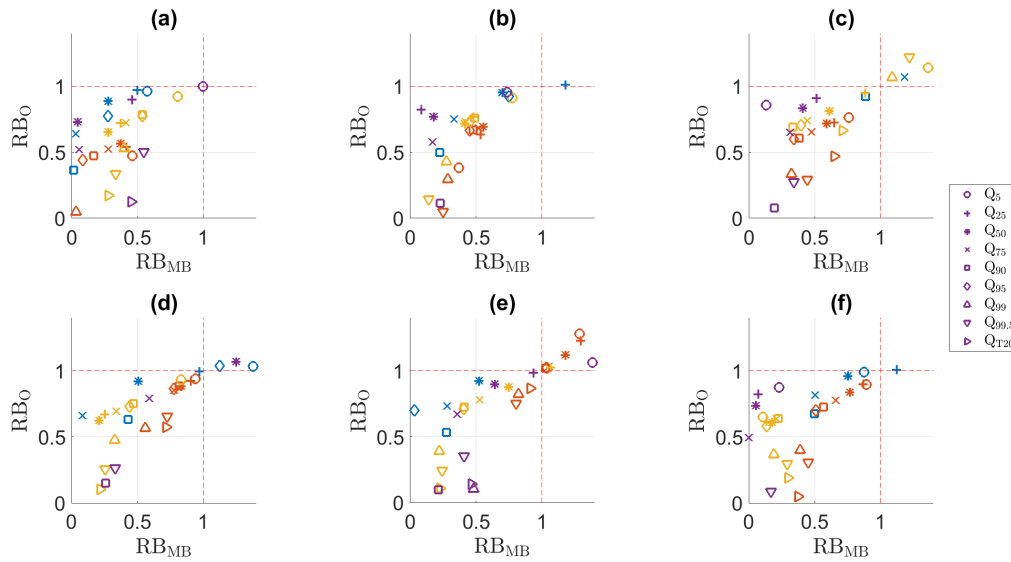
1005 Returning to the discharge, it might be interesting to discuss whether or not the adjustment of E is truly needed. On the  
one hand, this variable is the most affected by bias nonstationarity. On the other hand, discharge is far less influenced by this  
variable than by P or temporal structure. The discharge has been calculated for this setting with raw E, the result of which is  
shown in Fig. ???. The results depend on the method: for QDM and mQDM, raw E data slightly exacerbate the results, while  
for dOTC the percentiles are all improved. Only for MRQNBC and MBCn, the results are highly dependent on The impact of  
1010 bias nonstationarity varies between winter (blue) and summer (Fig. 10, yellow). In winter, the impact is more clearly visible on  
the higher percentiles:  $Q_{99}$ ,  $Q_{99.5}$  and  $Q_{10001000T20}$  are all well adjusted in the considered percentile. For MBCn, calibration  
period by QDM, mQDM, dOTC and  $R^2D^2$ , but are much worse adjusted in the validation period. In summer, the 5th percentile  
and the 20-year return period value (with  $RB_O \leq 0$ ) are improved, whereas the 95th and 99th percentile  $RB_O$  and  $RB_{MB}$  values  
are deteriorated. For MRQNBC, the results are opposite: the 5th percentile  $RB_O$  and  $RB_{MB}$  values are deteriorated, and the  
1015 95th and 99th values are improved impact seems to be similar for all the percentiles.

The results for raw E seem to imply that, on the one hand and depending on the bias-adjusting method used, a well-considered  
choice of variables to adjust can give optimal results. On the other hand, the results demonstrate once more that the univariate  
methods are far more robust than the multivariate methods. Although the  $RB_O$  and  $RB_{t10001000MB}$  values are slightly  
deteriorated for QDM and mQDM in comparison with the discharge based on adjusted E, all values, and especially the values  
1020 for the highest percentiles, still indicate a good bias adjustment. In general, an assessment like this can be done for the other  
types of impact studies discussed above, so that the influence of adjusting bias nonstationary variables can be better understood.

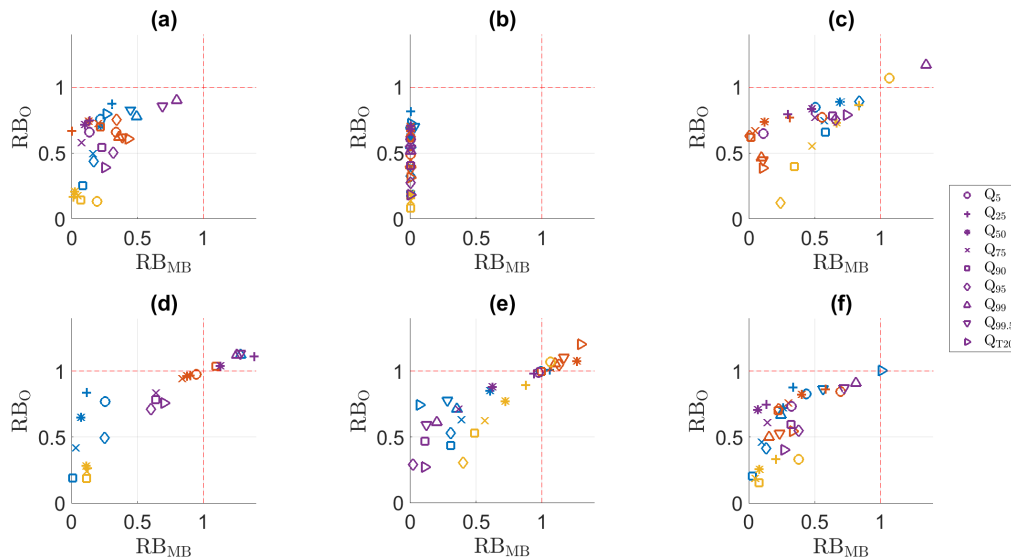
## 5 Discussion and conclusions

The goal of this paper was to assess how five-six bias-adjusting methods handle a climate change context with possible bias  
1025 nonstationarity. Three-Four of the methods were multivariate bias-adjusting methods: MRQNBC, MBCn and dOTC, dOTC and  
 $R^2D^2$ . The two other the bias-adjusting methods-ones were univariate: one was a traditional bias-adjusting method (QDM),

$RB_{MB}$  versus  $RB_0$  for the discharge percentiles and the 20-year return period value, calculated with raw evaporation. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.



**Figure 10.**  $RB_{MB}$  versus  $RB_0$  for the discharge percentiles and the 20-year return period value in the validation period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.



**Figure 11.**  $RB_{MB}$  versus  $RB_0$  for the discharge percentiles and the 20-year return period value in the calibration period. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC, (f)  $R^2D^2$ . Winter: blue, spring: ochre, summer: yellow, autumn: purple.



while the other was almost the same method, but modified according to the delta change paradigm (mQDM). These univariate methods were used as a baseline ~~to compare the multivariate bias-adjusting methods with for comparison~~. The climate change context, using 1970-1989 as calibration time period and 1998-2017 as validation time period, allowed us to calculate the change in bias between the periods, or the extent of bias ~~stationarity nonstationarity~~, using the R index. ~~All methods were calculated and The results of all methods were~~ compared using different indices, for which the residual biases relative to the observations and model bias were calculated. ~~Although the study was limited in spatial scale and climate models used, this yielded some results that could be valuable starting points for future research.~~

The calculated R index values ~~differed depending on the variable and variable index under consideration, but~~ generally demonstrated that the bias of some of these indices is not stationary under climate change conditions. ~~These changes could in some cases,~~ ~~although the extent of bias nonstationarity depended on the variable and index under consideration.~~ The bias nonstationarity could be clearly linked to the poor performance of bias-adjusting methods ~~, such as for for precipitation, and to some extent for temperature and potential evaporation.~~ For both precipitation and evaporation, it could be observed that the nonstationarity propagated through the rainfall-runoff model used for impact assessment, and that the propagation was different for these variables.

In the context of nonstationarity, it is important to discuss how well the methods performed. Some observations could be made. First, the ~~lower percentiles of T or the middle percentiles of E.~~ The performance was often poorer for the multivariate bias-adjusting methods, which corroborates the conclusions of Guo et al. (2020) that bias nonstationarity influences the performance of multivariate univariate bias-adjusting methods. ~~Although these methods have been developed during the last few years as a means to better adjust the biases, it seems that their more complex calculations make them more vulnerable to bias nonstationarity. Thus, the methods are relatively robust.~~ Although there always is an impact when bias nonstationarity is present, the univariate bias-adjusting methods ~~, computationally less complex and not taking (potentially changing) correlations into account,~~ seem to be more robust. ~~Although effective difference in climate change impact is weakened by the hydrological model we used, the univariate still perform best when considering the PSS values, i.e. the full PDF.~~ However, the methods are specifically designed to alter the marginal distributions. As already discussed in Section 4.5, it was pointed out by Zscheischler et al. (2019) that the multivariate bias-adjusting methods were made with other principal goals, such as spatial and dependence adjustment. As it is not assessed in this study, we cannot comment on the spatial adjustment. Nonetheless, the study by François et al. (2020) illustrated that the multivariate bias-adjusting methods still perform best. ~~Studying other types of climate change impacts, the effect of bias nonstationarity could possible be even larger than discussed here.~~ can be very informative and robust for spatial adjustment. Concerning the dependence adjustment, it was shown in Section 4.5 that the multivariate methods all perform well for the area and model chain studied here. Second, while QDM and mQDM seem to respond similarly, it should be taken into account that mQDM is designed to have a perfect fit in the calibration period. However, the poorer performance of mQDM for the precipitation occurrence indices is an indication that assuming that the temporal structure of the past can be used for the future might be dangerous, as Johnson and Sharma (2011) and Kerkhoff et al. (2014) already mentioned. Given that mQDM performed worse for two time periods separated by 10 years only, it is unlikely that it is safe to use this method, or other delta change-based methods, for impact assessments targeting the end of the 21st century that depend on the temporal structure of

time series. Yet, for some other indices, especially the correlation, mQDM performed better. Consequently, the exact choice should depend on the goals of the end user. Third, the methods with seasonal components do not always perform similarly. MRQNBC is able to address seasonal effects, but its performance varies strongly depending on the variable. Even in the situation where the univariate methods perform well, MRQNBC sometimes performed much worse, such as for temperature in autumn or in winter (Fig 4, panel (d), respectively purple and blue). Although these three observations can be made, it is impossible to fully discuss the method performance based on the set-up considered. The most important cause is the seasonality of the bias nonstationarity: while the bias nonstationarity shows clear differences between the seasons, some of the multivariate bias-adjusting methods are not yet equipped to handle seasonality. When there are large seasonal differences for the variables, for example for E and T, this causes a relatively poor performance in the calibration period, and a similar poor performance in the validation period. It is thus unclear whether the poor seasonal performance obfuscates the effect of nonstationarity, or if the similar performance is a sign of robustness. An earlier study (Guo et al., 2020) indicates the former, but this could also be location- and method-dependent. Hence, the set-up does not allow to clearly discern between the various categories of multivariate bias-adjustment, such as the ‘marginal/dependence’ or ‘all-in-one’ categories. To fully address the question on performance under bias nonstationarity, a better seasonal performance for the multivariate bias-adjusting methods seems crucial. However, not only seasonal differences in bias nonstationarity should be acknowledged: for variables other than P, T or E, or for other regions, bias nonstationarity might be better discernible on a monthly timescale, on a yearly timescale, or even on longer timescales. Only a few multivariate bias-adjusting methods specifically address multiple timescales, such as MRQNBC (Mehrotra and Sharma, 2016), or more recently, ‘Multivariate Frequency Bias Correction’ (MFBC) (Nguyen et al., 2018) or ‘3DBC’ (Mehrotra and Sharma, 2019). Yet, the varying performance of MRQNBC shows that the implementation of the seasonality can have a large impact. As such, the question about seasonality is not easy to answer.

The validation results could only be obtained by analysing and comparing a broad combination of indices. Considering only the mean or other standard statistics would have hidden many of the results seen. For example, in contrast to the results for the mean, the inclusion of both high and low extremes highlighted some problems with bias nonstationarity for some variables. As such, this study does not contradict earlier studies such as Maraun (2012), where the mean-based biases were found to be rather stable. ~~Even a broader set of indices, such as the ETCCDI indices, was not enough to clearly discern between the methods. As such, Thus,~~ we repeat the advice by Maraun and Widmann (2018a) to use indices not directly affected by bias-adjusting methods and to analyse the user needs before deciding upon the bias adjustment validation method. An important limitation is that we only used one GCM-RCM-combination. Using a model ensemble ~~will~~ would be more informative, but could hide a single model’s poor performance. On the other hand, similar assessments could also be used to discard poor-performing models (~~expanding upon methods such as those used in e.g. Brunner et al. (2019) or Tokarska et al. (2020)~~), based on the R index (also suggested by Maurer et al. (2013)) or the remaining bias after adjustment. ~~However, the used indices can still be improved.~~ Although the R index provides a lot of insight into the bias nonstationarity, it has been shown to over- or underestimate the effect of bias nonstationarity depending on the size and sometimes even the sign of the original bias. Other criteria also exist, such as the ‘signal-to-noise ratio’ (SNR) used by Hui et al. (2020). The different criteria or indices should be compared and maybe new tools are needed, so that the issue of bias nonstationarity can be more thoroughly explored.

The results for the multivariate bias-adjusting methods assessed here are in line with François et al. (2020), especially for the problematic adjustment of occurrence. François et al. (2020) consequently state that the different multivariate bias-adjusting methods are based on different assumptions, and thus, the To have a better view of how these results should be interpreted for impacts and compound events, the perspective of the end user should make well-grounded choices on the method used. This also became clear in our assessment. However, François et al. (2020) did not study the effect of climate change and bias (non)stationarity and instead focused on model trend preservation, or trend nonstationarity. The results presented and discussed here, such as the contrasting results of MRQNBC and dOTC, imply that whether trend preservation was the focus of a method or not, can be considered (Maraun et al., 2015; Maraun and Widmann, 2018b). We used discharge as an example, using the relatively simple PDM. Even for this model, it could be observed that bias nonstationarity can propagate in multiple ways. The influence of the nonstationarity in precipitation was most clear in summer and winter. As precipitation is the driving variable for the PDM, even the limited nonstationarity, mostly in the precipitation extremes, had an influence on the bias adjustment discharge simulation, as could be seen for the discharge in winter and summer (Fig. 11, respectively blue and yellow). In contrast, the nonstationarity in evaporation propagated much less. However, it is yet unclear how trend nonstationarity and bias nonstationarity influence each other and how the most appropriate methods can be discerned had an effect on the full PDF in spring, as could be observed from the PSS value for discharge (Table 6). In spring, no nonstationarity could be observed for precipitation, which allowed the influence of evaporation to be larger, although it has been suggested to use trend-preserving methods whenever we can assume the models to correctly simulate the atmospheric processes (Maraun, 2016)

Although critical of their use, the results of this paper do not imply that multivariate bias-adjusting methods are not helpful. Many of the methods developed during the past few years can also be used for spatial bias adjustment, in which case the locations can be used as extra variables (see theoretically has a smaller influence than precipitation on the discharge. The different propagation of bias nonstationarity, observed here for the extremes versus the full PDF, can be important considering that bias adjustment can be applied for many different types of impact assessment. However, the assessment in this study is relatively simple. For other impact studies, the results may vary considerably. For example, forest fires (a typical compound event, discussed in a bias adjustment context in e.g. Vrac (2018)) Yang et al. (2015), Cannon (2018), Zscheischler et al. (2019) depend more heavily on T and E to simulate fire weather conditions. Besides such compound events, other types of application can use a wide variety of variables and, hence, the bias nonstationarity may differ. In all of these studies, the propagation of bias nonstationarity will depend on the timescales considered in the impact assessment, the timescales on which nonstationarity is present, the variables considered and the spatial scale. Although this last aspect is limited in this study, it can be assumed that if bias nonstationarity is present in one grid cell, it will also be present in neighbouring grid cells with similar climatic conditions. A similar set-up has not been tested here, but the study by François et al. (2020) has proven the multivariate bias-adjusting methods to be very informative and robust for spatial adjustment: the spatial characteristics that influence local weather the most, such as orography.

Nonetheless To conclude, the results discussed in this paper indicate that many methods, and especially the multivariate bias-adjusting methods, fail in handling climate change and its resulting bias nonstationarity correctly bias nonstationarity can

1135 have an important influence on the bias adjustment and the propagation of biases in impact models. Depending on the extent of nonstationarity (spatial, temporal and the variables affected), such propagation should be taken into account far more when studying future impacts. As authors have mentioned before (Ehret et al., 2012; Maraun, 2016; Nahar et al., 2017), this foremost implies that climate models have to become better at modelling the future: we need to be able to trust them as fully as possible. As long as this is not the case, bias adjustment methods have to be developed that are more robust and that are able to help us assessing the future correctly. As such, the issue of seasonality as raised here is very important. Yet, impact assessment cannot wait for new methods to be developed and/or tested: we need to prepare ourselves for the future as soon as possible. ~~As was shown here~~For now, we can state ~~for the current generation of methods that the fewer assumptions and calculations a method needs, the more robust it is when used in a climate change context~~that for a robust bias adjustment under bias-nonstationary conditions, accounting for seasonality is crucial. Given this statement, we advise to use univariate bias-adjusting methods, until it becomes more clear how it can be ensured that multivariate methods certainly perform well ~~in a climate change context~~under bias nonstationarity.

1145 *Code and data availability.* The code for the computations is publicly available at <https://doi.org/10.5281/zenodo.4247518> (Hydro-Climate Extremes Lab – Ghent University, 2020). All methods were implemented in Matlab, except R2D2, for which the R package "R2D2" was used (Vrac and Thao, 2020a). The RCA4 data are downloaded and are available from the Earth System Grid Federation data repository. The local observations were obtained from RMI in Belgium, and cannot be shared with third parties.

*Author contributions.* JVDV, BDB and NV designed the experiments. JVDV developed the code and performed the calculations. JVDV prepared the manuscript with contributions from MD, BDB and NV. All co-authors contributed to the interpretation of the results.

1150 *Competing interests.* The authors declare that they have no conflict of interests.

*Acknowledgements.* J. Van de Velde would like to thank Y. Robin and R. Mehrotra for some helpful discussion on the use of respectively dOTC and MRQNBC. The authors are grateful to the RMI for allowing the use of 117-year Uccle dataset. This work was funded by FWO, grant number G.0039.18N. We would also like to thank Bastien François and one anonymous referee for their constructive comments and Carlo De Michele for editing.

## 1155 **6 Appendix A**

~~Observed values, and biases for the raw and adjusted climate simulations: (continued)~~ **Index Observed value QDM mQDM MBCn MRQNBC dOTC** P<sub>5</sub> (mm) 0.00 0.00 0.0000 0.0000 0.00 -0.40 0 P<sub>25</sub> (mm) 0.00 0.08 0.0000 0.0000 0.00 0.00 0.42

1160  $P_{50}$  (mm) 0.10 1.01 0.05 0.10 0.05 0.27 0.83  $P_{75}$  (mm) 2.70 1.83 -0.18 -0.17 -0.18 -0.84 0.76  $P_{90}$  (mm) 7.40 1.99 -0.26 -0.30  
 -0.26 -1.36 -0.23  $P_{95}$  (mm) 11.42 2.38 -0.61 -0.60 -0.61 -1.44 -0.65  $P_{99}$  (mm) 21.80 2.36 -1.86 -1.73 -1.86 1.38 -1.55  $P_{99.5}$   
 (mm) 29.09 1.56 -4.20 -3.97 -4.20 0.02 -4.02  $T_5$  (°C) 0.40 -0.31 -0.70 -1.43 -0.70 -0.63 -0.94  $T_{25}$  (°C) 6.30 -0.08 -0.68 -1.55  
 -0.68 -3.00 -0.73  $T_{50}$  (°C) 11.40 -0.40 -0.81 -0.88 -0.81 -3.24 -0.70  $T_{75}$  (°C) 16.10 -0.70 -0.46 -0.09 -0.46 -1.07 -0.37  $T_{90}$  (°C)  
 19.40 -1.07 -0.35 0.31 -0.35 1.02 0.12  $T_{95}$  (°C) 21.30 -1.17 -0.01 0.37 -0.01 2.59 0.25  $T_{99}$  (°C) 24.95 -1.85 -0.16 -0.75 -0.16  
 5.95 0.46  $T_{99.5}$  (°C) 25.90 -1.80 0.16 -1.39 0.16 7.01 0.77  $E_5$  (mm) 0.00 0.20 0.00 0.00 0.00 -0.04 0  $E_{25}$  (mm) 0.52 0.15 -0.09  
 -0.10 -0.09 -0.16 -0.52  $E_{50}$  (mm) 1.42 0.05 -0.27 -0.28 -0.27 -0.38 -0.77  $E_{75}$  (mm) 2.69 -0.02 -0.34 -0.35 -0.34 -0.58 -0.65  
 1165  $E_{90}$  (mm) 3.65 0.10 -0.27 -0.27 -0.27 -0.47 -0.18  $E_{95}$  (mm) 4.21 0.15 -0.30 -0.28 -0.30 -0.41 0.125  $E_{99}$  (mm) 5.02 0.21 -0.16  
 -0.13 -0.16 -0.17 0.69  $E_{99.5}$  (mm) 5.24 0.27 -0.10 -0.02 -0.10 -0.07 1.03  $corr_{P,E}$  (-) -0.18 -0.04 -0.06 0.02 0.19 0.17 0.57  
 $corr_{P,T}$  (-) -0.16 0.18 0.14 0.09 0.16 0.16 0.04  $corr_{E,T}$  (-) 0.82 -0.02 -0.03 -0.09 -0.84 0.02 -0.45  $crosseorr_{P,E,0}$  (-) 0.30 0.06  
 -0.04 -0.02 0.05 -0.03 0.12  $crosseorr_{P,T,0}$  (-) 0.24 0.19 0.08 -0.01 0.10 0.05 0.11  $crosseorr_{E,T,0}$  (-) 0.36 0.14778813 0.05 0.01  
 0.02 -0.03 0.06  $crosseorr_{P,E,T}$  (-) 0.38 0.126718335 0.02 -0.01 0.02 -0.04 0.08  $crosseorr_{P,T,T}$  (-) 0.93 -0.001694362 -0.02  
 1170 -0.05 -0.29 -0.02 -0.21  $crosseorr_{E,T,T}$  (-) 0.91 0.007385905 -0.01 -0.04 -0.27 -0.01 -0.24  $P_{P00}$  (-) 0.65 -0.10 -0.00 -0.02 -0.17  
 0.00 -0.37  $P_{P10}$  (-) 0.32 -0.15 0.00 -0.05 0.16 -0.13 -0.07  $N_{dry}$  (-) 3470.00 -1466.00 0.00 -373.00 0.00 -923.00 -1604.20  $P_{tag1}$   
 (-) 0.33 0.11 0.02 0.07 -0.12 0.08 0.05  $Q_5$  (m<sup>3</sup>/s) 2.30 0.92 -0.32 -0.18 0.82 -0.40 1.50  $Q_{25}$  (m<sup>3</sup>/s) 3.36 1.45 0.02 -0.12 0.29  
 -0.23 1.50  $Q_{50}$  (m<sup>3</sup>/s) 4.39 1.53 0.08 0.02 -0.20 -0.42 1.33  $Q_{75}$  (m<sup>3</sup>/s) 5.72 2.52 -0.08 -0.03 -0.72 -0.81 1.51  $Q_{90}$  (m<sup>3</sup>/s) 7.83  
 4.76 -0.36 -0.07 -1.66 -1.37 2.12  $Q_{95}$  (m<sup>3</sup>/s) 10.09 9.22 -1.00 -0.33 -2.78 -1.78 2.94  $Q_{99}$  (m<sup>3</sup>/s) 18.71 18.58 -1.65 -0.78 -3.21  
 1175 4.77 5.77  $Q_{99.5}$  (m<sup>3</sup>/s) 23.90 19.70 0.84 -1.77 -0.57 13.81 6.61  $Q_{T20}$  (m<sup>3</sup>/s) 48.69 54.61 8.36 -3.41 -10.40 52.45 25.03

## References

- Addor, N. and Fischer, E. M.: The influence of natural variability and interpolation errors on bias characterization in RCM simulations, *Journal of Geophysical Research: Atmospheres*, 120, 10–180, <https://doi.org/10.1002/2014JD022824>, 2015.
- Addor, N. and Seibert, J.: Bias correction for hydrological impact studies – beyond the daily perspective, *Hydrological Processes*, 28, 4823–4828, <https://doi.org/10.1002/hyp.10238>, 2014.
- Argüeso, D., Evans, J. P., and Fita, L.: Precipitation bias correction of very high resolution regional climate models, *Hydrology and Earth System Sciences*, 17, 4379, <https://doi.org/10.5194/hess-17-4379-2013>, 2013.
- Bárdossy, A. and Pegram, G.: Multiscale spatial recorrelation of RCM precipitation to produce unbiased climate change scenarios over large areas and small, *Water Resources Research*, 48, W09 502, <https://doi.org/10.1029/2011WR011524>, 2012.
- 1185 Beck, H. E., Zimmermann, N. E., McVicar, T. R., Vergopolan, N., Berg, A., and Wood, E. F.: Present and future Köppen-Geiger climate classification maps at 1-km resolution, *Scientific data*, 5, 180 214, <https://doi.org/10.1038/sdata.2018.214>, 2018.
- Bellprat, O., Kotlarski, S., Lüthi, D., and Schär, C.: Physical constraints for temperature biases in climate models, *Geophysical Research Letters*, 40, 4042–4047, <https://doi.org/10.1002/grl.50737>, 2013.
- Berg, P., Feldmann, H., and Panitz, H.-J.: Bias correction of high resolution regional climate model data, *Journal of Hydrology*, 448, 80–92, <https://doi.org/10.1016/j.jhydrol.2012.04.026>, 2012.
- 1190 Brunner, L., Lorenz, R., Zumwald, M., and Knutti, R.: Quantifying uncertainty in European climate projections using combined performance-independence weighting, *Environmental Research Letters*, 14, 124 010, <https://doi.org/10.1088/1748-9326/ab492f>, 2019.
- Buser, C. M., Künsch, H. R., Lüthi, D., Wild, M., and Schär, C.: Bayesian multi-model projection of climate: bias assumptions and interannual variability, *Climate Dynamics*, 33, 849–868, <https://doi.org/10.1007/s00382-009-0588-6>, 2009.
- 1195 Cabus, P.: River flow prediction through rainfall–runoff modelling with a probability-distributed model (PDM) in Flanders, Belgium, *Agricultural Water Management*, 95, 859–868, <https://doi.org/10.1016/j.agwat.2008.02.013>, 2008.
- Cannon, A. J.: Multivariate bias correction of climate model output: Matching marginal distributions and intervariable dependence structure, *Journal of Climate*, 29, 7045–7064, <https://doi.org/10.1175/JCLI-D-15-0679.1>, 2016.
- Cannon, A. J.: Multivariate quantile mapping bias correction: an N-dimensional probability density function transform for climate model simulations of multiple variables, *Climate Dynamics*, 50, 31–49, <https://doi.org/10.1007/s00382-017-3580-6>, 2018.
- 1200 Cannon, A. J., Sobie, S. R., and 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, 6938–6959, <https://doi.org/10.1175/JCLI-D-14-00754.1>, 2015.
- Chen, J., Brissette, F. P., and Lucas-Picher, P.: Assessing the limits of bias-correcting climate model outputs for climate change impact studies, *Journal of Geophysical Research: Atmospheres*, 120, 1123–1136, <https://doi.org/10.1002/2014JD022635>, 2015.
- 1205 Christensen, J. H., Boberg, F., Christensen, O. B., and Lucas-Picher, P.: On the need for bias correction of regional climate change projections of temperature and precipitation, *Geophysical Research Letters*, 35, L20 709, <https://doi.org/10.1029/2008GL035694>, 2008.
- Clark, M., Gangopadhyay, S., Hay, L., Rajagopalan, B., and Wilby, R.: The Schaake shuffle: A method for reconstructing space–time variability in forecasted precipitation and temperature fields, *Journal of Hydrometeorology*, 5, 243–262, [https://doi.org/10.1175/1525-7541\(2004\)005<0243:TSSAMF>2.0.CO;2](https://doi.org/10.1175/1525-7541(2004)005<0243:TSSAMF>2.0.CO;2), 2004.
- 1210 Cuturi, M.: Sinkhorn distances: Lightspeed computation of optimal transport, in: *Advances in Neural Information Processing Systems*, pp. 2292–2300, 2013.

- De Jongh, I. L. M., Verhoest, N. E. C., and De Troch, F. P.: Analysis of a 105-year time series of precipitation observed at Uccle, Belgium, *International Journal of Climatology*, 26, 2023–2039, <https://doi.org/10.1002/joc.1352>, 2006.
- 1215 Dekens, L., Parey, S., Grandjacques, M., and Dacunha-Castelle, D.: Multivariate distribution correction of climate model outputs: A generalization of quantile mapping approaches, *Environmetrics*, 28, e2454, <https://doi.org/10.1002/env.2454>, 2017.
- Demarée, G. R.: The centennial recording raingauge of the Uccle Plateau: Its history, its data and its applications, *Houille Blanche*, 4, 95–102, 2003.
- Derbyshire, J.: The siren call of probability: Dangers associated with using probability for consideration of the future, *Futures*, 88, 43–54, <https://doi.org/10.1016/j.futures.2017.03.011>, 2017.
- 1220 Di Luca, A., de Elía, R., and Laprise, R.: Challenges in the quest for added value of regional climate dynamical downscaling, *Current Climate Change Reports*, 1, 10–21, <https://doi.org/10.1007/s40641-015-0003-9>, 2015.
- Eberhart, R. and Kennedy, J.: A new optimizer using Particle Swarm Theory, in: *Proceedings of the Sixth International Symposium on Micro Machine and Human Science*, pp. 39–43, IEEE, 1995.
- Ehret, U., Zehe, E., Wulfmeyer, V., Warrach-Sagi, K., and Liebert, J.: HESS Opinions" Should we apply bias correction to global and regional climate model data?", *Hydrology and Earth System Sciences*, 16, 3391–3404, <https://doi.org/10.5194/hess-16-3391-2012>, 2012.
- 1225 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geoscientific Model Development*, 9, 1937–1958, <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.
- Fosser, G., Kendon, E. J., Stephenson, D., and Tucker, S.: Convection-permitting models offer promise of more certain extreme rainfall projections, *Geophysical Research Letters*, p. e2020GL088151, <https://doi.org/10.1029/2020GL088151>, 2020.
- 1230 François, B., Vrac, M., Cannon, A. J., Robin, Y., and Allard, D.: Multivariate bias corrections of climate simulations: Which benefits for which losses?, *Earth System Dynamics*, 2020, 1–41, <https://doi.org/10.5194/esd-11-537-2020>, 2020.
- Galmarini, S., Cannon, A., Ceglar, A., Christensen, O., Dentener, F. J., de Noblet-Ducoudré, N., Doblas-Reyes, F. J., and Vrac, M.: Adjusting climate model bias for agricultural impact assessment: How to cut the mustard, *Climate Services*, 13, 65–69, <https://doi.org/10.1016/j.cliser.2019.01.004>, 2019.
- 1235 Gudmundsson, L., Bremnes, J. B., Haugen, J. E., and Engen-Skaugen, T.: Downscaling RCM precipitation to the station scale using statistical transformations—a comparison of methods, *Hydrology and Earth System Sciences*, 16, 3383–3390, <https://doi.org/10.5194/hess-16-3383-2012>, 2012.
- Guo, Q., Chen, J., Zhang, X. J., Xu, C.-Y., and Chen, H.: Impacts of using state-of-the-art multivariate bias correction methods on hydrological modeling over North America, *Water Resources Research*, 56, e2019WR026659, <https://doi.org/10.1029/2019WR026659>, 2020.
- 1240 Gutiérrez, J. M., Maraun, D., Widmann, M., Huth, R., Hertig, E., Benestad, R., Roessler, O., Wibig, J., Wilcke, R., Kotlarski, S., Martín, D. S., Herrera, S., Bedia, J., Casanueva, A., Manzanar, R., Iturbide, M., Vrac, M., Dubrovsky, M., Ribalaygua, J., Pórtoles, J., Rätty, O., Räisänen, J., Hingray, B., Raynaud, D., Casado, M. J., Ramos, P., Zerenner, T., Turco, M., Bosshard, T., Štěpánek, P., Bartholy, J., Pongracz, R., Keller, D. E., Fischer, A. M., Cardoso, R. M., Soares, P. M. M., Czernecki, B., and Pagé, C.: An intercomparison of a large ensemble of statistical downscaling methods over Europe: results from the VALUE perfect predictor cross-validation experiment, *International Journal of Climatology*, 39, 3750–3785, <https://doi.org/10.1002/joc.5462>, 2019.
- 1245 Gutjahr, O. and Heinemann, G.: Comparing precipitation bias correction methods for high-resolution regional climate simulations using COSMO-CLM, *Theoretical and Applied Climatology*, 114, 511–529, <https://doi.org/10.1007/s00704-013-0834-z>, 2013.

- Gutowski, William J., J., Decker, S. G., Donavon, R. A., Pan, Z., Arritt, R. W., and Takle, E. S.: Temporal–spatial scales of  
1250 observed and simulated precipitation in central US climate, *Journal of Climate*, 16, 3841–3847, [https://doi.org/10.1175/1520-0442\(2003\)016<3841:TSSOAS>2.0.CO;2](https://doi.org/10.1175/1520-0442(2003)016<3841:TSSOAS>2.0.CO;2), 2003.
- Haerter, J., Hagemann, S., Moseley, C., and Piani, C.: Climate model bias correction and the role of timescales, *Hydrology and Earth System Sciences*, 15, 1065–1073, <https://doi.org/10.5194/hess-15-1065-2011>, 2011.
- Hagemann, S., Chen, C., Haerter, J. O., Heinke, J., Gerten, D., and Piani, C.: Impact of a statistical bias correction on the pro-  
1255 jected hydrological changes obtained from three GCMs and two hydrology models, *Journal of Hydrometeorology*, 12, 556–578, <https://doi.org/10.1175/2011JHM1336.1>, 2011.
- Hakala, K., Addor, N., and Seibert, J.: Hydrological modeling to evaluate climate model simulations and their bias correction, *Journal of Hydrometeorology*, 19, 1321–1337, <https://doi.org/10.1175/JHM-D-17-0189.1>, 2018.
- Hay, L. E. and Clark, M. P.: Use of statistically and dynamically downscaled atmospheric model output for hydrologic simulations in three  
1260 mountainous basins in the western United States, *Journal of Hydrology*, 282, 56–75, [https://doi.org/10.1016/S0022-1694\(03\)00252-X](https://doi.org/10.1016/S0022-1694(03)00252-X), 2003.
- Helsen, S., van Lipzig, N. P. M., Demuzere, M., Vanden Broucke, S., Caluwaerts, S., De Cruz, L., De Troch, R., Hamdi, R., Termonia, P., Van Schaeybroeck, B., and Wouters, H.: Consistent scale-dependency of future increases in hoclimate models, *Climate Dynamics*, 54, 1–14, <https://doi.org/10.1007/s00382-019-05056-w>, 2019.
- 1265 Hempel, S., Frieler, K., Warszawski, L., Schewe, J., and Piontek, F.: A trend-preserving bias correction—the ISI-MIP approach, *Earth System Dynamics*, 4, 219–236, <https://doi.org/10.5194/esd-4-219-2013>, 2013.
- Hewitson, B. C., Daron, J., Crane, R. G., Zermoglio, M. F., and Jack, C.: Interrogating empirical-statistical downscaling, *Climatic change*, 122, 539–554, <https://doi.org/10.1007/s10584-013-1021-z>, 2014.
- Higham, N. J.: Computing a nearest symmetric positive semidefinite matrix, *Linear Algebra and its Applications*, 103, 103–118,  
1270 [https://doi.org/10.1016/0024-3795\(88\)90223-6](https://doi.org/10.1016/0024-3795(88)90223-6), 1988.
- Ho, C. K., Stephenson, D. B., Collins, M., Ferro, C. A. T., and Brown, S. J.: Calibration strategies: a source of additional uncertainty in climate change projections, *Bulletin of the American Meteorological Society*, 93, 21–26, <https://doi.org/10.1175/2011BAMS3110.1>, 2012.
- Hui, Y., Chen, J., Xu, C.-Y., Xiong, L., and Chen, H.: Bias nonstationarity of global climate model outputs: The role of internal climate variability and climate model sensitivity, *International Journal of Climatology*, 39, 2278–2294, <https://doi.org/10.1002/joc.5950>, 2019.
- 1275 Hui, Y., Xu, Y., Chen, J., Xu, C.-Y., and Chen, H.: Impacts of bias nonstationarity of climate model outputs on hydrological simulations, *Hydrology Research*, 51, 925–941, <https://doi.org/10.2166/nh.2020.254>, 2020.
- Hydro-Climate Extremes Lab – Ghent University: h-cel/ImpactofBiasNonstationarity: Impact of bias nonstationarity: calculations, <https://doi.org/10.5281/zenodo.4247518>, 2020.
- Ines, A. V. M. and Hansen, J. W.: Bias correction of daily GCM rainfall for crop simulation studies, *Agricultural and Forest Meteorology*,  
1280 138, 44–53, <https://doi.org/10.1016/j.agrformet.2006.03.009>, 2006.
- IPCC: Managing the Risks of Extreme Events and Disasters to Advance Climate Change Adaptation: Special Report of the Intergovernmental Panel on Climate Change, Cambridge University Press, 2012.
- IPCC: Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University Press, 2013.
- 1285 Ivanov, M. A., Luterbacher, J., and Kotlarski, S.: Climate model biases and modification of the climate change signal by intensity-dependent bias correction, *Journal of Climate*, 31, 6591–6610, <https://doi.org/10.1175/JCLI-D-17-0765.1>, 2018.



- Jacob, D., Petersen, J., Eggert, B., Alias, A., Christensen, O. B., Bouwer, L. M., Braun, A., Colette, A., Déqué, M., Georgievski, G., Georgopoulou, E., Gobiet, A., Menut, L., Nikulin, G., Haensler, A., Hempelmann, N., Jones, C., Keuler, K., Kovats, S., Kröner, N., Kotlarski, S., Kriegsmann, A., Martin, E., van Meijgaard, E., Moseley, C., Pfeifer, S., Preuschmann, S., Radermacher, C., Radtke, K., Rechid, D., Rounsevell, M., Samuelsson, P., Somot, S., Soussana, J.-F., Teichmann, C., Valentini, R., Vautard, R., Weber, B., and Yiou, P.: EURO-CORDEX: new high-resolution climate change projections for European impact research, *Regional Environmental Change*, 14, 563–578, <https://doi.org/10.1007/s10113-013-0499-2>, 2014.
- 1290 Johnson, F. and Sharma, A.: Accounting for interannual variability: A comparison of options for water resources climate change impact assessments, *Water Resources Research*, 47, W04 508, <https://doi.org/10.1029/2010WR009272>, 2011.
- 1295 Johnson, F. and Sharma, A.: A nesting model for bias correction of variability at multiple time scales in general circulation model precipitation simulations, *Water Resources Research*, 48, W01 504, <https://doi.org/10.1029/2011WR010464>, 2012.
- Kendon, E. J., Ban, N., Roberts, N. M., Fowler, H. J., Roberts, M. J., Chan, S. C., Evans, J. P., Fosser, G., and Wilkinson, J. M.: Do convection-permitting regional climate models improve projections of future precipitation change?, *Bulletin of the American Meteorological Society*, 98, 79–93, <https://doi.org/10.1175/BAMS-D-15-0004.1>, 2017.
- 1300 Kerkhoff, C., Künsch, H. R., and Schär, C.: Assessment of bias assumptions for climate models, *Journal of Climate*, 27, 6799–6818, <https://doi.org/10.1175/JCLI-D-13-00716.1>, 2014.
- Klemeš, V.: Operational testing of hydrological simulation models, *Hydrological Sciences Journal*, 31, 13–24, <https://doi.org/10.1080/02626668609491024>, 1986.
- Knol, D. L. and ten Berge, J. M. F.: Least-squares approximation of an improper correlation matrix by a proper one, *Psychometrika*, 54, 53–61, <https://doi.org/10.1007/BF02294448>, 1989.
- 1305 Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet, A., Goergen, K., Jacob, D., Lüthi, D., van Meijgaard, E., Nikulin, G., Schär, C., Teichmann, C., Vautard, R., Warrach-Sagi, K., and Wulfmeyer, V.: Regional climate modeling on European scales: a joint standard evaluation of the EURO-CORDEX RCM ensemble, *Geoscientific Model Development*, 7, 1297–1333, <https://doi.org/10.5194/gmd-7-1297-2014>, 2014.
- 1310 Lenderink, G., Buishand, A., and Van Deursen, W.: Estimates of future discharges of the river Rhine using two scenario methodologies: direct versus delta approach, *Hydrology and Earth System Sciences*, 11, 1145–1159, <https://doi.org/10.5194/hess-11-1145-2007>, 2007.
- Li, C., Sinha, E., Horton, D. E., Diffenbaugh, N. S., and Michalak, A. M.: Joint bias correction of temperature and precipitation in climate model simulations, *Journal of Geophysical Research: Atmospheres*, 119, 13–153, <https://doi.org/10.1002/2014JD022514>, 2014.
- 1315 Li, H., Sheffield, J., and 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 101, <https://doi.org/10.1029/2009JD012882>, 2010.
- Lorenz, E. N.: Atmospheric predictability as revealed by naturally occurring analogues, *Journal of Atmospheric Sciences*, 26, 636–646, [https://doi.org/10.1175/1520-0469\(1969\)26<636:APARBN>2.0.CO;2](https://doi.org/10.1175/1520-0469(1969)26<636:APARBN>2.0.CO;2), 1969.
- Maraun, D.: Nonstationarities of regional climate model biases in European seasonal mean temperature and precipitation sums, *Geophysical Research Letters*, 39, <https://doi.org/https://doi.org/10.1029/2012GL051210>, 2012.
- 1320 Maraun, D.: Bias correcting climate change simulations—a critical review, *Current Climate Change Reports*, 2, 211–220, <https://doi.org/10.1007/s40641-016-0050-x>, 2016.
- Maraun, D. and Widmann, M.: Cross-validation of bias-corrected climate simulations is misleading, *Hydrology and Earth System Sciences*, 22, 4867–4873, <https://doi.org/10.5194/hess-22-4867-2018>, 2018a.

- 1325 Maraun, D. and Widmann, M.: Statistical Downscaling and Bias Correction for Climate Research, Cambridge University Press, <https://doi.org/10.1017/9781107588783>, 2018b.
- Maraun, D., Widmann, M., Gutiérrez, J. M., Kotlarski, S., Chandler, R. E., Hertig, E., Wibig, J., Huth, R., and Wilcke, R. A. I.: VALUE: A framework to validate downscaling approaches for climate change studies, *Earth's Future*, 3, 1–14, <https://doi.org/10.1002/2014EF000259>, 2015.
- 1330 Matalas, N. C.: Mathematical assessment of synthetic hydrology, *Water Resources Research*, 3, 937–945, <https://doi.org/10.1029/WR003i004p00937>, 1967.
- Maurer, E. P., Das, T., and Cayan, D. R.: Errors in climate model daily precipitation and temperature output: time invariance and implications for bias correction, *Hydrology and Earth System Sciences*, 17, 2147–2159, <https://doi.org/10.5194/hess-17-2147-2013>, 2013.
- Mehrotra, R. and Sharma, A.: An improved standardization procedure to remove systematic low frequency variability biases in GCM simulations, *Water Resources Research*, 48, W12 601, <https://doi.org/10.1029/2012WR012446>, 2012.
- 1335 Mehrotra, R. and Sharma, A.: Correcting for systematic biases in multiple raw GCM variables across a range of timescales, *Journal of Hydrology*, 520, 214–223, <https://doi.org/10.1016/j.jhydrol.2014.11.037>, 2015.
- Mehrotra, R. and Sharma, A.: A multivariate quantile-matching bias correction approach with auto- and cross-dependence across multiple time scales: Implications for downscaling, *Journal of Climate*, 29, 3519–3539, <https://doi.org/10.1175/jcli-d-15-0356.1>, 2016.
- 1340 Mehrotra, R. and Sharma, A.: A Resampling Approach for Correcting Systematic Spatiotemporal Biases for Multiple Variables in a Changing Climate, *Water Resources Research*, 55, 754–770, <https://doi.org/10.1029/2018WR023270>, 2019.
- Meyer, J., Kohn, I., Stahl, K., Hakala, K., Seibert, J., and Cannon, A. J.: Effects of univariate and multivariate bias correction on hydrological impact projections in alpine catchments, *Hydrology and Earth System Sciences*, 23, 1339–1354, <https://doi.org/10.5194/hess-23-1339-2019>, 2019.
- 1345 Michelangeli, P.-A., Vrac, M., and Loukos, H.: Probabilistic downscaling approaches: Application to wind cumulative distribution functions, *Geophysical Research Letters*, 36, L11 708, <https://doi.org/10.1029/2009GL038401>, 2009.
- Milly, P. C. D., Betancourt, J., Falkenmark, M., Hirsch, R. M., Kundzewicz, Z. W., Lettenmaier, D. P., and Stouffer, R. J.: Stationarity is dead: Whither water management?, *Science*, 319, 573–574, <https://doi.org/10.1126/science.1151915>, 2008.
- Moore, R. J.: The PDM rainfall-runoff model, *Hydrology and Earth System Sciences*, 11, 483–499, <https://doi.org/10.5194/hess-11-483-2007>, 2007.
- 1350 Nahar, J., Johnson, F., and Sharma, A.: Assessing the extent of non-stationary biases in GCMs, *Journal of Hydrology*, 549, 148–162, <https://doi.org/10.1016/j.jhydrol.2017.03.045>, 2017.
- Nelsen, R. B.: *An Introduction to Copulas*, 2nd, New York: Springer Science Business Media, 2006.
- Nguyen, H., Mehrotra, R., and Sharma, A.: Correcting for systematic biases in GCM simulations in the frequency domain, *Journal of Hydrology*, 538, 117–126, <https://doi.org/10.1016/j.jhydrol.2016.04.018>, 2016.
- 1355 Nguyen, H., Mehrotra, R., and Sharma, A.: Correcting systematic biases across multiple atmospheric variables in the frequency domain, *Climate Dynamics*, 52, 1283–1298, <https://doi.org/10.1007/s00382-018-4191-6>, 2018.
- Olsson, J., Berggren, K., Olofsson, M., and Viklander, M.: Applying climate model precipitation scenarios for urban hydrological assessment: A case study in Kalmar City, Sweden, *Atmospheric Research*, 92, 364–375, <https://doi.org/10.1016/j.atmosres.2009.01.015>, 2009.
- 1360 Panofsky, H. A., Brier, G. W., and Best, W. H.: *Some Application of Statistics to Meteorology*, Earth and Mineral Sciences Continuing Education, College of Earth and Mineral Sciences, Pennsylvania State University, 1958.

- Papalexiou, S. M. and Montanari, A.: Global and regional increase of precipitation extremes under global warming, *Water Resources Research*, 55, 4901–4914, <https://doi.org/10.1029/2018WR024067>, 2019.
- Penman, H. L.: Natural evaporation from open water, bare soil and grass, *Proc. R. Soc. Lond. A*, 193, 120–145, 1948.
- 1365 Perkins, S. E., Pitman, A. J., Holbrook, N. J., and McAneney, J.: Evaluation of the AR4 climate models' simulated daily maximum temperature, minimum temperature, and precipitation over Australia using probability density functions, *Journal of climate*, 20, 4356–4376, <https://doi.org/https://doi.org/10.1175/JCLI4253.1>, 2007.
- Peyré, G. and Cuturi, M.: *Computational Optimal Transport*, vol. 11, Now Publishers, <https://doi.org/10.1561/22000000073>, 2019.
- Pham, M. T.: Copula-based stochastic modelling of evapotranspiration time series conditioned on rainfall as design tool in water resources management, PhD thesis, Faculty of Biosciences Engineering, Ghent University, 2016.
- 1370 Pham, M. T., Vernieuwe, H., De Baets, B., and Verhoest, N. E. C.: A coupled stochastic rainfall–evapotranspiration model for hydrological impact analysis, *Hydrology and Earth System Sciences*, 22, 1263–1283, <https://doi.org/10.5194/hess-22-1263-2018>, 2018.
- Piani, C. and Haerter, J. O.: Two dimensional bias correction of temperature and precipitation copulas in climate models, *Geophysical Research Letters*, 39, <https://doi.org/10.1029/2012gl053839>, 2012.
- 1375 Piani, C., Haerter, J. O., and Coppola, E.: Statistical bias correction for daily precipitation in regional climate models over Europe, *Theoretical and Applied Climatology*, 99, 187–192, <https://doi.org/10.1007/s00704-009-0134-9>, 2010.
- Popke, D., Stevens, B., and Voigt, A.: Climate and climate change in a radiative-convective equilibrium version of ECHAM6, *Journal of Advances in Modeling Earth Systems*, 5, 1–14, <https://doi.org/10.1029/2012MS000191>, 2013.
- Prein, A. F., Langhans, W., Fosser, G., Ferrone, A., Ban, N., Goergen, K., Keller, M., Tölle, M., Gutjahr, O., Feser, F., Brisson, E., Kollet, S., 1380 Schmidli, J., Lipzig, N. P. M., and Leung, R.: A review on regional convection-permitting climate modeling: Demonstrations, prospects, and challenges, *Reviews of Geophysics*, 53, 323–361, <https://doi.org/10.1002/2014RG000475>, 2015.
- Rajczak, J., Kotlarski, S., and Schär, C.: Does quantile mapping of simulated precipitation correct for biases in transition probabilities and spell lengths?, *Journal of Climate*, 29, 1605–1615, <https://doi.org/10.1175/JCLI-D-15-0162.1>, 2016.
- Räty, O., Räisänen, J., and Ylhäisi, J. S.: Evaluation of delta change and bias correction methods for future daily precipitation: intermodel cross-validation using ENSEMBLES simulations, *Climate Dynamics*, 42, 2287–2303, <https://doi.org/10.1007/s00382-014-2130-8>, 2014.
- 1385 Räty, O., Räisänen, J., Bosshard, T., and Donnelly, C.: Intercomparison of univariate and joint bias correction methods in changing climate from a hydrological perspective, *Climate*, 6, 33, <https://doi.org/10.3390/cli6020033>, 2018.
- Reiter, P., Gutjahr, O., Schefczyk, L., Heinemann, G., and Casper, M.: Does applying quantile mapping to subsamples improve the bias correction of daily precipitation?, *International Journal of Climatology*, 38, 1623–1633, <https://doi.org/10.1002/joc.5283>, 2018.
- 1390 Rizzo, M. L. and Székely, G. J.: Energy distance, *Wiley Interdisciplinary Reviews: Computational Statistics*, 8, 27–38, <https://doi.org/10.1002/wics.1375>, 2016.
- Robin, Y., Vrac, M., Naveau, P., and Yiou, P.: Multivariate stochastic bias corrections with optimal transport, *Hydrology and Earth System Sciences*, 23, 773–786, <https://doi.org/10.5194/hess-23-773-2019>, 2019.
- Rojas, R., Feyen, L., Dosio, A., and Bavera, D.: Improving pan-European hydrological simulation of extreme events through statistical bias correction of RCM-driven climate simulations, *Hydrology & Earth System Sciences*, 15, <https://doi.org/10.5194/hess-15-2599-2011>, 2011.
- 1395 Salas, J. D.: *Applied Modeling of Hydrologic Time Series*, Water Resources Publication, 1980.
- Schmidli, J., Frei, C., and Vidale, P. L.: Downscaling from GCM precipitation: a benchmark for dynamical and statistical downscaling methods, *International Journal of Climatology*, 26, 679–689, <https://doi.org/10.1002/joc.1287>, 2006.

- 1400 Schölzel, C. and Friederichs, P.: Multivariate non-normally distributed random variables in climate research-introduction to the copula approach, *Nonlinear Processes in Geophysics*, 15, 761–772, <https://doi.org/10.5194/npg-15-761-2008>, 2008.
- Shepherd, T. G.: Atmospheric circulation as a source of uncertainty in climate change projections, *Nature Geoscience*, 7, 703–708, <https://doi.org/10.1038/ngeo2253>, 2014.
- Sippel, S., Otto, F. E. L., Forkel, M., Allen, M. R., Guillod, B. P., Heimann, M., Reichstein, M., Seneviratne, S. I., Thonicke, K.,  
1405 and Mahecha, M. D.: A novel bias correction methodology for climate impact simulations, *Earth System Dynamics*, 7, 71–88, <https://doi.org/10.5194/esd-7-71-2016>, 2016.
- Srikanthan, R. and Pegram, G. G. S.: A nested multisite daily rainfall stochastic generation model, *Journal of Hydrology*, 371, 142–153, <https://doi.org/10.1016/j.jhydrol.2009.03.025>, 2009.
- Strandberg, G., Barring, L., Hansson, U., Jansson, C., Jones, C., Kjellström, E., Kupiainen, M., Nikulin, G., Samuelsson, P., and Ullerstig,  
1410 A.: CORDEX scenarios for Europe from the Rossby Centre regional climate model RCA4, Tech. rep., SMHI, 2015.
- Sunyer, M. A., Madsen, H., Rosbjerg, D., and Arnbjerg-Nielsen, K.: A Bayesian approach for uncertainty quantification of extreme precipitation projections including climate model interdependency and nonstationary bias, *Journal of Climate*, 27, 7113–7132, <https://doi.org/10.1175/JCLI-D-13-00589.1>, 2014.
- Switanek, M. B., Troch, P. A., Castro, C. L., Leuprecht, A., Chang, H. I., Mukherjee, R., and Demaria, E. M. C.: Scaled distribution mapping:  
1415 A bias correction method that preserves raw climate model projected changes, *Hydrology and Earth System Sciences*, 21, 2649–2666, <https://doi.org/10.5194/hess-21-2649-2017>, 2017.
- Teutschbein, C. and 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, <https://doi.org/10.1016/j.jhydrol.2012.05.052>, 2012.
- Teutschbein, C. and Seibert, J.: Is bias correction of regional climate model (RCM) simulations possible for non-stationary conditions?,  
1420 *Hydrology and Earth System Sciences*, 17, 5061–5077, <https://doi.org/10.5194/hess-17-5061-2013>, 2013, 2013.
- Themeßl, M. J., Gobiet, A., and Leuprecht, A.: Empirical-statistical downscaling and error correction of daily precipitation from regional climate models, *International Journal of Climatology*, 31, 1530–1544, <https://doi.org/10.1002/joc.2168>, 2011.
- Themeßl, M. J., Gobiet, A., and Heinrich, G.: Empirical-statistical downscaling and error correction of regional climate models and its impact on the climate change signal, *Climatic Change*, 112, 449–468, <https://doi.org/10.1007/s10584-011-0224-4>, 2012.
- 1425 Tokarska, K. B., Stolpe, M. B., Sippel, S., Fischer, E. M., Smith, C. J., Lehner, F., and Knutti, R.: Past warming trend constrains future warming in CMIP6 models, *Science Advances*, 6, eaaz9549, <https://doi.org/10.1126/sciadv.aaz9549>, 2020.
- Van de Velde, J., De Baets, B., Demuzere, M., and Verhoest, N. E. C.: Comparison of occurrence-bias-adjusting methods for hydrological impact modelling, *Hydrology and Earth System Sciences Discussions*, 2020, 1–35, <https://doi.org/10.5194/hess-2020-83>, 2020.
- Van Schaeybroeck, B. and Vannitsem, S.: Assessment of calibration assumptions under strong climate changes, *Geophysical Research Letters*, 43, 1314–1322, <https://doi.org/10.1002/2016GL067721>, 2016.  
1430
- van Vuuren, D. P., Edmonds, J., Kainuma, M., Riahi, K., Thomson, A., Hibbard, K., Hurtt, G. C., Kram, T., Krey, V., Lamarque, J.-F., Masui, T., Meinshausen, M., Nakicenovic, N., Smith, S. J., and Rose, S. K.: The representative concentration pathways: an overview, *Climatic Change*, 109, 5, <https://doi.org/10.1007/s10584-011-0148-z>, 2011.
- Vandenbergh, S., Verhoest, N. E. C., Onof, C., and De Baets, B.: A comparative copula-based bivariate frequency analysis of observed and simulated storm events: A case study on Bartlett-Lewis modeled rainfall, *Water Resources Research*, 47, W07 529, 2011.  
1435
- Velázquez, J. A., Troin, M., Caya, D., and Brissette, F.: Evaluating the time-invariance hypothesis of climate model bias correction: implications for hydrological impact studies, *Journal of Hydrometeorology*, 16, 2013–2026, <https://doi.org/10.1175/JHM-D-14-0159.1>, 2015.

- Verhoest, N. E. C., Troch, P. A., and De Troch, F. P.: On the applicability of Bartlett–Lewis rectangular pulses models in the modeling of design storms at a point, *Journal of Hydrology*, 202, 108–120, [https://doi.org/10.1016/S0022-1694\(97\)00060-7](https://doi.org/10.1016/S0022-1694(97)00060-7), 1997.
- 1440 Verstraeten, G., Poesen, J., Demarée, G. R., and Salles, C.: Long-term (105 years) variability in rain erosivity as derived from 10-min rainfall depth data for Ukkel (Brussels, Belgium): Implications for assessing soil erosion rates, *Journal of Geophysical Research*, 111, D22 109, <https://doi.org/10.1029/2006jd007169>, 2006.
- Villani, C.: *Optimal transport: old and new*, vol. 338, Springer Science & Business Media, 2008.
- Vrac, M.: Multivariate bias adjustment of high-dimensional climate simulations: the Rank Resampling for Distributions and Dependences (R2D2) bias correction, *Hydrology and Earth System Sciences*, 22, 3175, <https://doi.org/10.5194/hess-22-3175-2018>, 2018.
- 1445 Vrac, M. and Friederichs, P.: Multivariate—intervariable, spatial, and temporal—bias correction, *Journal of Climate*, 28, 218–237, <https://doi.org/10.1175/JCLI-D-14-00059.1>, 2015.
- Vrac, M. and Thao, S.: R package R2D2, <https://doi.org/10.5281/ZENODO.4021981>, 2020a.
- Vrac, M. and Thao, S.: R2D2 v2.0: accounting for temporal dependences in multivariate bias correction via analogue rank resampling, *Geoscientific Model Development*, 13, 5367–5387, <https://doi.org/10.5194/gmd-13-5367-2020>, 2020b.
- 1450 Vrac, M., Noël, T., and Vautard, R.: Bias correction of precipitation through Singularity Stochastic Removal: Because occurrences matter, *Journal of Geophysical Research: Atmospheres*, 121, 5237–5258, <https://doi.org/10.1002/2015JD024511>, 2016.
- Wang, L. and Chen, W.: Equiratio cumulative distribution function matching as an improvement to the equidistant approach in bias correction of precipitation, *Atmospheric Science Letters*, 15, 1–6, <https://doi.org/10.1002/asl2.454>, 2014.
- 1455 Wang, Y., Sivandran, G., and Bielicki, J. M.: The stationarity of two statistical downscaling methods for precipitation under different choices of cross-validation periods, *International Journal of Climatology*, 38, e330–e348, <https://doi.org/10.1002/joc.5375>, 2018.
- Wilcke, R. A. I., Mendlik, T., and Gobiet, A.: Multi-variable error correction of regional climate models, *Climatic Change*, 120, 871–887, <https://doi.org/10.1007/s10584-013-0845-x>, 2013.
- Willems, P.: Revision of urban drainage design rules after assessment of climate change impacts on precipitation extremes at Uccle, Belgium, *Journal of Hydrology*, 496, 166–177, <https://doi.org/10.1016/j.jhydrol.2013.05.037>, 2013.
- 1460 Willems, P. and Vrac, M.: Statistical precipitation downscaling for small-scale hydrological impact investigations of climate change, *Journal of Hydrology*, 402, 193–205, <https://doi.org/10.1016/j.jhydrol.2011.02.030>, 2011.
- Yang, W., Gardelin, M., Olsson, J., and Bosshard, T.: Multi-variable bias correction: application of forest fire risk in present and future climate in Sweden, *Natural Hazards and Earth System Sciences*, 15, 2037–2057, <https://doi.org/10.5194/nhess-15-2037-2015>, 2015.
- 1465 Zhang, X., Alexander, L., Hegerl, G. C., Jones, P., Tank, A. K., Peterson, T. C., Trewin, B., and Zwiers, F. W.: Indices for monitoring changes in extremes based on daily temperature and precipitation data, *Wiley Interdisciplinary Reviews: Climate Change*, 2, 851–870, <https://doi.org/10.1002/wcc.147>, 2011.
- Zorita, E. and Von Storch, H.: The analog method as a simple statistical downscaling technique: Comparison with more complicated methods, *Journal of climate*, 12, 2474–2489, [https://doi.org/10.1175/1520-0442\(1999\)012<2474:TAMAAS>2.0.CO;2](https://doi.org/10.1175/1520-0442(1999)012<2474:TAMAAS>2.0.CO;2), 1999.
- 1470 Zscheischler, J., Westra, S., Hurk, B. J. J. M., Seneviratne, S. I., Ward, P. J., Pitman, A., AghaKouchak, A., Bresch, D. N., Leonard, M., Wahl, T., and Zhang, X.: Future climate risk from compound events, *Nature Climate Change*, p. 1, <https://doi.org/10.1038/s41558-018-0156-3>, 2018.
- Zscheischler, J., Fischer, E. M., and Lange, S.: The effect of univariate bias adjustment on multivariate hazard estimates, *Earth System Dynamics*, 10, 31–43, <https://doi.org/10.5194/esd-10-31-2019>, 2019.

- 1475 Zscheischler, J., Martius, O., Westra, S., Bevacqua, E., Raymond, C., Horton, R. M., van den Hurk, B., AghaKouchak, A., Jézéquel, A., Mahecha, M. D., Maraun, D., Ramos, A. M., Ridder, N. N., Thiery, W., and Vignotto, E.: A typology of compound weather and climate events, *Nature Reviews Earth & Environment*, <https://doi.org/10.1038/s43017-020-0060-z>, 2020.