

**Observational uncertainty and regional climate model  
evaluation: A pan-European perspective**

Journal:	<i>International Journal of Climatology</i>
Manuscript ID	JOC-17-0256.R1
Wiley - Manuscript type:	VALUE Special Issue
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Kotlarski, Sven; Federal Institute of Meteorology and Climatology MeteoSwiss, Climate</p> <p>Szabó, Péter; Hungarian Meteorological Service, Climate and Ambient Air</p> <p>Herrera, Sixto; Universidad de Cantabria, Meteorology Group, Dpto. de Matemática Aplicada y Ciencias de la Computación</p> <p>Räty, Olle; Univ Helsinki, Finland, Department</p> <p>Keuler, Klaus; Brandenburg University of Technology, Chair Environmental Meteorology</p> <p>Soares, Pedro; Instituto Dom Luiz, Universidade de Lisboa, DEGGE</p> <p>Cardoso, Rita; Instituto Dom Luiz, Universidade de Lisboa, DEGGE</p> <p>Bosshard, Thomas; Swedish Meteorological and Hydrological Institute, Hydrology Department</p> <p>Pagé, Christian; UMR CNRS 5318 CECI – CERFACS, Climate modeling and Global change</p> <p>Boberg, Fredrik; Danish Meteorological Institute, Danish Climate Centre</p> <p>Gutiérrez, José; National Research Council (CSIC), Meteorology Group, Instituto de Física de Cantabria</p> <p>Isotta, Francesco; Federal Institute of Meteorology and Climatology MeteoSwiss, Climate</p> <p>Jaczewski, Adam; National Research Institute, Institute of Meteorology and Water Management</p> <p>Kreienkamp, Frank; Deutscher Wetterdienst, Klima- und Umweltberatung</p> <p>Liniger, Mark Andrea; Federal Office of Meteorology and Climatology MeteoSwiss, Climate</p> <p>Lussana, Cristian; Norwegian Meteorological Institute, Observation and Climate</p> <p>Pianko-Kluczyńska, Krystyna; National Research Institute, Institute of Meteorology and Water Management</p>
Keywords:	RCM evaluation, observations, uncertainty, Europe, CORDEX
Country Keywords:	Germany, France, Switzerland, Spain, Sweden

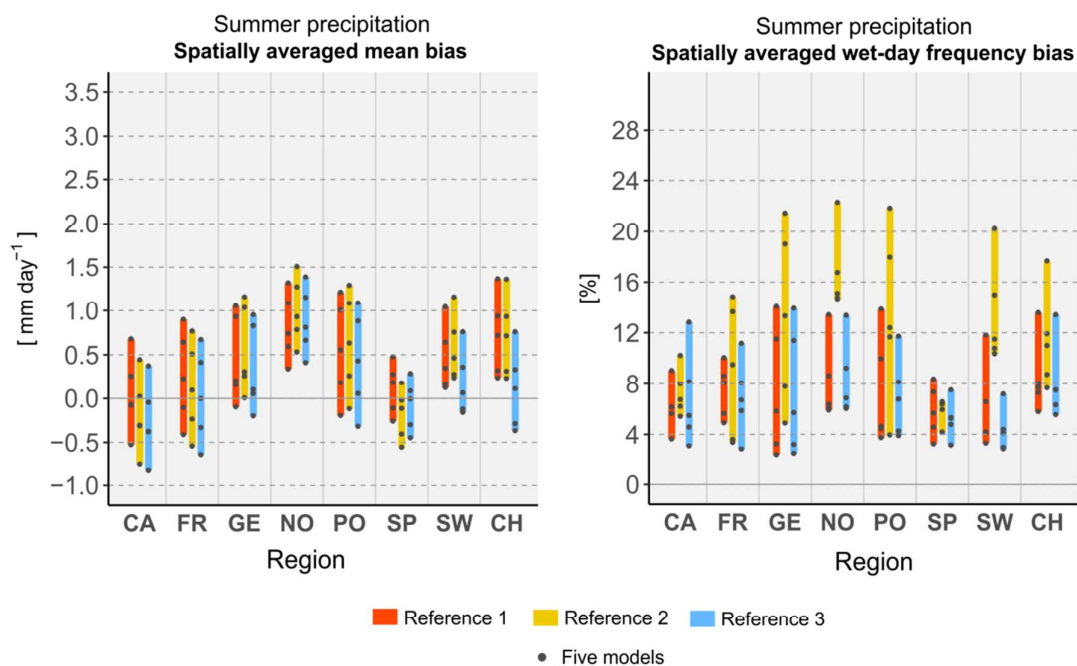
SCHOLARONE™  
Manuscripts

Peer Review Only

# Observational uncertainty and regional climate model evaluation: A pan-European perspective

Sven Kotlarski\*, Péter Szabó, Sixto Herrera, Olle Räty, Klaus Keuler, Pedro M. Soares, Rita M. Cardoso, Thomas Bosshard, Christian Pagé, Fredrik Boberg, José M. Gutiérrez, Francesco A. Isotta, Adam Jaczewski, Frank Kreienkamp, Mark A. Liniger, Cristian Lussana, Krystyna Pianko-Kluczyńska

Five state-of-the-art reanalysis-driven regional climate model experiments are evaluated against three different observational reference datasets for two variables (temperature and precipitation) and for eight sub-regions of the European continent. Overall, we find the influence of observational uncertainty to be smaller than model uncertainty. For individual regions and seasons, however, model evaluation can considerably depend on the chosen reference and final model ranks can be strongly influenced.



# Observational uncertainty and regional climate model evaluation: A pan-European perspective

Sven Kotlarski<sup>1</sup>, Péter Szabó<sup>2</sup>, Sixto Herrera<sup>3</sup>, Olle Rätty<sup>4</sup>, Klaus Keuler<sup>5</sup>, Pedro M. Soares<sup>6</sup>, Rita M. Cardoso<sup>6</sup>, Thomas Bosshard<sup>7</sup>, Christian Pagé<sup>8</sup>, Fredrik Boberg<sup>9</sup>, José M. Gutiérrez<sup>10</sup>, Francesco A. Isotta<sup>1</sup>, Adam Jaczewski<sup>11</sup>, Frank Kreienkamp<sup>12</sup>, Mark A. Liniger<sup>1</sup>, Cristian Lussana<sup>13</sup>, Krystyna Pianko-Kluczyńska<sup>11</sup>

<sup>1</sup>Federal Office of Meteorology and Climatology MeteoSwiss, Switzerland

<sup>2</sup>Hungarian Meteorological Service, Hungary

<sup>3</sup>Meteorology Group, Dpto. de Matemática Aplicada y Ciencias de la Computación. Universidad de Cantabria, Spain

<sup>4</sup>University of Helsinki, Finland

<sup>5</sup>Brandenburg University of Technology, Germany

<sup>6</sup>Instituto Dom Luiz, Faculdade de Ciências, Universidade de Lisboa, Portugal

<sup>7</sup>Swedish Meteorological and Hydrological Institute, Sweden

<sup>8</sup>UMR CNRS 5318 CECI – CERFACS, France

<sup>9</sup>Danish Meteorological Institute, Denmark

<sup>10</sup>Meteorology Group, Instituto de Física de Cantabria (CSIC-Univ. de Cantabria), Spain

<sup>11</sup>Institute of Meteorology and Water Management – National Research Institute, Poland

<sup>12</sup>Deutscher Wetterdienst, Germany

<sup>13</sup>Norwegian Meteorological Institute, Norway

## Abstract

The influence of uncertainties in gridded observational reference data on regional climate model (RCM) evaluation is quantified on a pan-European scale. Three different reference datasets are considered: the coarse-resolved E-OBS dataset, a compilation of regional high-resolution gridded products (HR) and the European-scale MESAN reanalysis. Five high-resolution ERA-Interim driven RCM experiments of the EURO-CORDEX initiative are evaluated against each of these references over eight European sub-regions and considering a range of performance metrics for mean daily temperature and daily precipitation. The spatial scale of the evaluation is 0.22°, i.e. the grid spacing of the coarsest dataset in the exercise (E-OBS).

While the three reference grids agree on the overall mean climatology, differences can be pronounced over individual regions. These differences partly translate into RCM evaluation uncertainty. Still, for most cases observational uncertainty is smaller than RCM uncertainty. For individual sub-regions and performance metrics, however, observational uncertainty can dominate. This is especially true for precipitation and for metrics targeting the wet-day frequency, the pattern correlation and the distributional similarity. In some cases also the spatially averaged mean bias can be considerably affected.

An illustrative ranking exercise highlights the overall effect of observational uncertainty on RCM ranking. Over individual sub-domains, the choice of a specific reference can modify RCM ranks by up to four levels (out of five RCMs). For most cases, however, RCM ranks are stable irrespective of the reference. These results provide a two-fold picture: model uncertainty dominates for most regions and for most performance metrics considered, and observational uncertainty plays a minor role. For individual cases, however, observational uncertainty can be pronounced and needs to be definitely taken into account. Results can to some extent also depend on the treatment of potential precipitation undercatch in the observational reference.

## 46 **Keywords**

47 RCM evaluation, observations, uncertainty, Europe, CORDEX

## 48 **1. Introduction**

49 The existence and availability of reliable high-quality observational data is essential for climate  
50 monitoring. It is furthermore the basis for the development, evaluation and application of both  
51 physically-based and statistical weather and climate models. This includes downscaling approaches  
52 that translate large-scale atmospheric features into higher-resolved and even point-scale information  
53 (e.g., Fowler et al., 2007). Observations are already used during model development, but also model  
54 calibration and initialization often heavily rely on an existing observational reference (e.g., Bellprat et  
55 al., 2012). As such, the quality of any model-derived weather or climate product can be expected to  
56 depend on the quality of the underlying observations. The same is true for model evaluation  
57 exercises that assess and inter-compare the performance of one or several modelling systems by  
58 comparison against observation-based records (e.g., Christensen et al., 2010, Kotlarski et al., 2014).  
59 Consequently, uncertainties in the observational reference directly translate into uncertainties of  
60 model evaluation results.

61 Observational uncertainties themselves can be large and originate from multiple sources. Already  
62 raw observations are likely to suffer from inaccuracies due to residual non-climatic influences  
63 (Hartmann et al., 2013, Hegerl et al., 2001, McMillan et al., 2012). Such influences include  
64 malfunctions and error margins of measurement devices and, in case of long-term records,  
65 replacements of the device, relocations of the measurement site or physical changes of the  
66 surrounding landscape. For the case of precipitation, site measurements are furthermore subject to  
67 systematic biases due to the local deformation of the wind field by the gauge and wetting and  
68 evaporation losses. This systematic undercatch is pronounced for windy conditions and for snowfall  
69 and can result in an important underestimation of true precipitation sums (e.g., Adam and  
70 Lettenmaier, 2003, Cheval et al., 2010, Frei et al., 2003, Groisman and Legates, 1994, Sevruk, 1985,  
71 Wolff et al., 2015). Some of the mentioned inaccuracies can be reduced by postprocessing the raw  
72 measurement records, e.g. by applying data homogenization procedures (Begert et al., 2005) or a  
73 dedicated precipitation undercatch correction (Richter, 1995). Additionally, representativity issues  
74 arise for point measurements, i.e. the question to what extent a point record reflects conditions for a  
75 larger area, for instance the mean conditions over a climate model grid box obtained through  
76 averaging all subgrid variabilities in space (e.g., Osborn and Hulme, 1997).

77 To avoid the latter complication, climate model evaluation wherever possible relies on gridded  
78 reference datasets that are obtained by a spatial analysis and interpolation of point measurements  
79 onto a regular grid yielding area-representative grid cell mean values. Additionally, gridded remote  
80 sensing products and model-derived reanalyses are used. In any case, the gridding procedure itself  
81 involves assumptions and uncertainties with corresponding effects on the final product. For gridded  
82 datasets obtained by spatial interpolation of point measurements problems arise especially in  
83 regions with sparse data coverage, complex topography and for variables with a high spatio-temporal  
84 climatic variability (e.g., Wagner et al., 2007). Spatial variance, for instance, is mostly underestimated  
85 by gridded products (Beguería et al., 2016) and trends can be affected by a temporally changing  
86 network density (e.g., Frei, 2014, Hofstra et al., 2009). Sampling issues due to random natural climate  
87 variability, i.e. the fact that the observed record is only one possible realization of the analysis  
88 period's climate, can introduce further uncertainties (e.g., Addor and Fischer, 2015, Mahlstein et al.,  
89 2015).

90 In summary, any available observation-based record is unlikely to reflect the true state of  
91 atmospheric quantities but only some approximation of it. A number of studies exist that quantify  
92 the related observational uncertainty by comparing several observation-based reference datasets for

93 specific variables and regions (e.g., the recent works by Awange et al., 2016, Berg et al., 2016, Dunn  
94 et al., 2014, Gbambie et al., 2017, Gervais et al., 2014, Herold et al., 2016, Hofstra et al., 2009, Isotta  
95 et al., 2015, Kyselý and Plavcová, 2010, Palazzi et al., 2013, Rauthe et al., 2013, Schneider et al., 2014;  
96 Tanarhte et al., 2012). In evaluation exercises these shortcomings of the reference inevitably  
97 influence the performance assessment of climate models and introduce uncertainties in the  
98 evaluation results. Previous works have addressed this issue by employing multiple reference data  
99 sources for global and regional climate model (GCM, RCM) evaluation (Addor and Fischer, 2015,  
100 Bellprat et al., 2012, Brienen et al., 2016, Bucchignani et al., 2016, Casanueva et al., 2013, Cheneka  
101 et al., 2016, Davin et al., 2016, Di Luca et al., 2012, Gómez-Navarro et al., 2012, Haslinger et al., 2013,  
102 Kotlarski et al., 2005, Kotlarski et al., 2012, Maraun et al., 2012, Prein and Gobiet, 2017, Ring et al.,  
103 2016, Sunyer et al., 2013). Besides quantifying the influence of observational uncertainty on  
104 individual model performance scores, two of these studies (Gómez-Navarro et al., 2012 and Sunyer  
105 et al., 2013) also explicitly address the modification of model ranks when changing the observational  
106 reference.

107 Most of the mentioned works consider geographic domains of limited extent only, such as individual  
108 river catchments or countries, and focus on precipitation. At this point, we refrain from listing the  
109 individual results but note that (1) even in regions covered by dense observational networks  
110 observational uncertainty can be large and can be comparable to RCM uncertainty (measured by the  
111 spread between individual RCM experiments) and that (2) observational uncertainty can have the  
112 potential to influence the outcome of climate model weighting and ranking exercises. Among the  
113 mentioned works a particularly relevant study is the one by Prein and Gobiet (2017) who, focusing on  
114 precipitation, inter-compared a large number of gridded observational datasets over parts of the  
115 European continent and used this observational ensemble to evaluate state-of-the-art RCM  
116 experiments. They found that observational uncertainty can be of similar magnitude as RCM biases,  
117 particularly in regions of low station density and for high temporal and spatial resolution statistics.

118 In the present work we build upon and complement these previous studies by quantifying  
119 observational uncertainty on a pan-European scale not only for precipitation but also for  
120 temperature and by assessing its influence on RCM evaluation in a well-defined performance  
121 assessment framework. We explicitly include an illustrative model ranking exercise and relate  
122 observational spread to RCM spread. Our main objective is to illustrate the influence of observational  
123 uncertainty on RCM evaluation and RCM ranking for different European sub-regions, for two  
124 variables and for a range of performance scores reflecting different model bias characteristics.

## 125 **2. Data and Methods**

### 126 **2.1 Observational Reference Data**

127 To sample observational uncertainty we employ three observational reference grids that are  
128 available (1) for both mean temperature and precipitation, (2) at a daily resolution, (3) for the  
129 common 18-year long evaluation period 1989-2006, and (4) at a grid spacing comparable to or higher  
130 than the current RCM resolution for multi-decadal climate projections. Note that the latter criterion  
131 does not necessarily imply a higher effective resolution of the observational datasets compared to  
132 the RCMs. Depending on the underlying network density the effective resolution of the data could be  
133 considerably lower than the nominal grid spacing (e.g., Beguería et al., 2016, Isotta et al., 2015, Prein  
134 and Gobiet, 2017) . The three observational reference grids represent an “ensemble of opportunity”,  
135 i.e. we consider datasets that are readily available, that fulfil the above-mentioned criteria and that  
136 include the evaluation of climate models in their intended range of application. We hence accept  
137 inter-dependencies of the three datasets that could arise, for instance, from the use of the same  
138 station series for gridding or calibration purposes or from similar gridding concepts. In particular, we  
139 combine reference datasets that result from an explicit gridding procedure of observations with a

140 reanalysis-based product. We also do not intend to provide final explanations for differences among  
141 the three reference datasets. This would imply a much more detailed analysis of the influence of the  
142 gridding process and of different network densities on the final gridded product and would go  
143 beyond the scope of the present work. These aspects are covered by the accompanying study of  
144 Herrera et al. (2017). Furthermore, note that we use the term *observations* for results from both  
145 gridding processes and reanalysis procedures. This contrasts with other, more direct definitions of  
146 *observations* based on actual station data or remote sensing results. We hence do not explicitly  
147 differentiate between *observational uncertainty* and *gridding uncertainty* and use the former term to  
148 capture both.

#### 149 2.1.1 E-OBS

150 The gridded E-OBS dataset (Haylock et al., 2008; version 15) covers the entire European land surface  
151 and is based on the ECA&D (European Climate Assessment and Dataset) station data plus more than  
152 2000 further stations from additional archives. We here use the daily temperature and precipitation  
153 grids of the rotated 0.22° version (approx. 25 km grid spacing). For several years E-OBS has now been  
154 a standard reference for RCM evaluation over the European continent. Known deficiencies of E-OBS  
155 relate to remaining inhomogeneities in the station series and to the dataset's quality in regions of  
156 sparse station density. The latter particularly affects the representation of daily extremes (e.g.,  
157 Bellprat et al., 2012; Herrera et al., 2012; Hofstra et al., 2009, Hofstra et al., 2010, Lenderink, 2010,  
158 Maraun et al., 2012) and the effective spatial resolution which is presumably lower than the nominal  
159 0.22° grid spacing (e.g., Hanel and Buishand, 2011; Kyselý and Plavcová, 2010). The systematic  
160 undercatch of rain gauges (e.g., Sevruk, 1986) has not been corrected for, i.e., E-OBS likely  
161 underestimates true precipitation sums.

#### 162 2.1.2 National High-Resolution Grids (HR)

163 Our second observational reference is a compilation of national/regional high-resolution  
164 temperature and precipitation grids that are available for parts of the European continent only (Fig.  
165 1). This dataset has been assembled within the COST Action VALUE (Maraun et al., 2015). It covers  
166 modified sets of regions and datasets compared to the recent work of Prein and Gobiet (2017),  
167 including one additional country (Poland), an updated version of the Norwegian and the German  
168 dataset and the consideration of Switzerland only instead of the entire Alps, employing a different  
169 high-resolution observational grid. In overlapping boarder regions covered by two national datasets  
170 only one of them has been considered<sup>1</sup>. In the following a brief description of each dataset is  
171 provided. Except for the Swedish product, none of the precipitation grids explicitly accounts for the  
172 systematic undercatch of rain gauges.

173 **Spain (SP):** For peninsular Spain and the Balearic Islands an improved 3-dimensional areal  
174 representative version (AA-3D) of the Spain02 gridded dataset at 0.22° grid spacing on a rotated grid  
175 is used (Herrera et al. 2012; 2016). Spain02 is based on a very dense and quality-controlled station  
176 network consisting of 2756 and 237 stations for precipitation and temperature, respectively. The  
177 interpolation and gridding procedure is the same as applied for E-OBS.

178 **Poland (PO):** The AA-3D methodology used for the Spanish grid was extended to build an  
179 observational grid for Poland based on a quality-controlled observational station dataset provided by  
180 the Institute of Meteorology and Water Management - National Research Institute, Center for  
181 Poland's Climate Monitoring; see Herrera et al. (2017) for further details. This dataset comprises 197  
182 stations for precipitation and 123 for temperature. Station data were homogenized prior to the

---

<sup>1</sup> In the following pairs of overlapping countries/regions the bold country/region has been considered: **NO**/SW,  
**SP**/FR, **CH**/GE, **FR**/GE, **CH**/FR, **PO**/CA, **GE**/PO. In case of the Carpathian dataset, which extends far into Poland,  
this means a substantial cut-off at its northern boundary.

183 gridding by applying the MASH v3.03 procedure (e.g., Szentimrey, 2013) to the daily data (Lakatos et  
184 al. 2013).

185 **France (FR):** The France national high-resolution analysis SAFRAN is available at an hourly time step  
186 and on a grid of 8 km spacing (Durand et al., 1993; Quintana-Seguí et al. 2008; Vidal et al., 2010). It is  
187 based on observations at more than 4000 sites collected by Météo-France as well as on operational  
188 Numerical Weather Prediction analyses along with some climatological data. It covers all water  
189 basins affecting Metropolitan France including Corsica. Prior to its use within the present work  
190 SAFRAN was conservatively interpolated to the rotated 0.11° EURO-CORDEX grid.

191 **Sweden (SW):** The daily gridded PTHBV dataset provides daily precipitation and temperature data at  
192 4 km grid spacing and covers Sweden plus some adjacent regions. The product is based on more than  
193 350 (800) stations for temperature (precipitation) and has been constructed by optimal interpolation  
194 with a climatological background field that accounts for wind-orography effects (Johansson and  
195 Chen, 2003). In the present work it is the only dataset that has been corrected to account for the  
196 systematic undercatch of rain gauges. The correction is based on gauge type, precipitation type (rain  
197 or snow), wind classification and exposure of the gauges (Berg et al., 2016).

198 **Germany (GE):** The high resolution daily gridded HYRAS dataset has been produced as part of the  
199 KLIWAS research programme (*Impacts of climate change on waterways and navigation – searching  
200 for options of adaptation*; [www.kliwas.de](http://www.kliwas.de)). It covers the period 1951 to 2006 and is available at 5 km  
201 grid spacing for all river catchments in Germany as well as adjacent river basins with drainage  
202 towards Germany (i.e. the entire Rhine, Danube and Elbe catchments). More detailed information  
203 about the dataset and its underlying station network, which consists of up to 1000 and 6200 stations  
204 for temperature and precipitation, respectively, is provided by Rauthe et al. (2013) and Frick et al.  
205 (2014).

206 **Carpathians (CA):** The CARPATCLIM gridded observational dataset (Lakatos et al., 2013) covers parts  
207 of 9 countries along the Carpathian Mountains and is based on raw station time series that were  
208 exchanged along the borders to ensure data homogeneity (temperature: 258 stations, precipitation:  
209 727 stations). Quality control and homogenization were carried out at daily resolution using the  
210 MASH software (Szentimrey, 2004). The MISH package (Szentimrey and Bihari, 2007) was employed  
211 for spatial interpolation. The publicly available CARPATCLIM dataset for 11 variables is provided at  
212 daily temporal resolution and 0.1° grid spacing for the period 1961-2010 ([www.carpatclim-eu.org](http://www.carpatclim-eu.org)).  
213 Note that, in contrast to the other national/regional grids, CARPATCLIM does not represent areal grid  
214 cell averages but point estimates for the grid cell centers.

215 **Norway (NO):** The gridded seNorge version 2 (seNorge2) dataset is based on two modified optimal  
216 interpolation schemes (Gandin, 1965), one for temperature and one for precipitation, in which the  
217 prior distribution is estimated from in-situ observations (Lussana et al., 2016, Uboldi et al., 2008). The  
218 input data used are original non-homogenized station series from the Norwegian Climate Database  
219 (480 and 920 stations on average for temperature and precipitation, respectively). Data for both  
220 variables are provided at daily resolution on a 1 km grid.

221 **Switzerland (CH):** For the region of Switzerland, the TabsD (temperature; MeteoSwiss, 2013a) and  
222 RhiresD (precipitation; MeteoSwiss, 2013b) datasets at 2 km grid spacing are used. Both datasets rely  
223 on a large but temporally varying number of station series (temperature: 93, precipitation: about  
224 520) and were produced accounting for the special requirements of interpolating station data in  
225 topographically complex terrain (e.g., Frei, 2014).

### 226 2.1.3 EURO4M MESAN

227 The European Reanalysis and Observations for Monitoring project (EURO4M) has produced several  
228 gridded datasets for Europe, among others a High Resolution Limited Area Model (HIRLAM)  
229 reanalysis at a grid spacing of 0.2° (approx. 22 km) using 3D-VAR data assimilation (Dahlgren and



230 Gustafsson, 2012). Several simulated surface fields – including near-surface air temperature and  
231 precipitation - have afterwards been downscaled with the MESAN system to a  $0.05^\circ$  grid (approx. 5  
232 km) using optimal interpolation techniques (Häggmark et al., 2000) and assimilating further surface  
233 observations. Depending on the region, the number of stations used for assimilation and  
234 interpolation is partly larger and partly smaller or comparable to E-OBS (see Fig. 1 in Prein and  
235 Gobiet, 2017). For precipitation, the surface observations assimilated in the MESAN downscaling step  
236 were not corrected for the measurement bias of rain gauges. Hence the final EURO4M MESAN  
237 precipitation product – although originating from simulated precipitation of the HIRLAM model – has  
238 to be assumed to be undercatch-affected.

## 239 2.2 RCM Data

240 The RCM simulations that are evaluated in the present work originate from the EURO-CORDEX  
241 initiative (Jacob et al., 2014) and have been carried out at a grid spacing of  $0.11^\circ$  on a rotated grid  
242 under the CORDEX simulation protocol. All experiments cover a full European domain (see Kotlarski  
243 et al., 2014) and were driven by the ERA-Interim reanalysis (Dee et al., 2011) at the lateral domain  
244 boundaries. We hence evaluate the so-called *perfect boundary experiments* instead of the GCM-  
245 driven historical control runs. Such an evaluation places a stronger focus on the downscaling  
246 performance itself as potentially strong biases in the GCM-derived boundary forcing are avoided. In  
247 total, five simulations are used (Table 1) that form a subset of those experiments considered in the  
248 EURO-CORDEX standard evaluation (Kotlarski et al., 2014). Note that two of the five RCMs employed  
249 (HIRHAM 5 and RACMO 2.2E) as well as the reanalysis-model (MESAN; see above) originate from the  
250 numerical weather prediction model HIRLAM and partly share the same code. Hence, their  
251 respective outputs cannot be considered to be fully independent of each other.

## 252 2.3 Analysis Domain and Analysis Grid

253 The analysis domain of the present work consists of the eight regions covered by HR (Fig. 1; Section  
254 2.1.2) and samples an important part of continental-scale climate variability in Europe. To enable a  
255 consistent comparison on a grid cell level the higher-resolved HR, MESAN and RCM data (including  
256 elevation) were conservatively aggregated to the rotated  $0.22^\circ$  E-OBS grid, i.e. to the coarsest grid  
257 considered in this work, prior to the analysis. This enables a grid-cell-by-grid-cell comparison and  
258 avoids the additional interpolation of E-OBS to the higher-resolved RCM grid. This procedure is also  
259 beneficial in case that the effective resolution of a certain dataset is smaller than its nominal grid  
260 spacing; spatial aggregation would then more accurately represent the effective resolution of the  
261 data. For temperature an additional elevation correction from the aggregated HR, MESAN and RCM  
262 elevation to the elevation of the corresponding E-OBS grid cell was carried out assuming a spatially  
263 and temporally uniform lapse rate of  $0.0065^\circ\text{C m}^{-1}$ .

## 264 2.4 Performance Metrics

265 The performance of the RCMs was evaluated on the common  $0.22^\circ$  analysis grid and separately for  
266 each of the eight sub-regions. Seven different metrics were chosen which describe different aspects  
267 of model performance. Five of these metrics were computed for both temperature and precipitation  
268 and one further metric was calculated for temperature or precipitation only, resulting in six metrics  
269 for each variable.

270 For each observational reference dataset the metrics were calculated for every climate model  $j$ ,  
271 season  $k$ , and analysis region  $r$ . For the sake of simplicity, those indices are omitted in the following.  
272 We define  $O_n$  and  $X_n$  to be daily observational and climate model data, respectively, at a particular  
273 grid point  $n$  within the analysis region  $r$  that contains a total of  $N$  grid points. Further, overbars

274 denote the temporal mean over all time steps in the analysis period that fall into the season  $k$ , two  
275 overbars denote temporal and spatial mean, and yearly seasonal means are denoted by the index  $y$ .

276 The performance of the climatological seasonal mean averaged over a sub-region was evaluated by  
277 the bias given as

$$278 \quad \mathbf{BIAS} = \frac{1}{N} \sum_{n=1}^N (\bar{X}_n - \bar{O}_n) \quad (\text{Eq. 1})$$

279 Moderate extremes at the upper end of the distribution were evaluated by the mean absolute error  
280 of the 99<sup>th</sup> percentile:

$$281 \quad \mathbf{MAE99} = \frac{1}{N} \sum_{n=1}^N |P^{99}(X_n) - P^{99}(O_n)| \quad (\text{Eq. 2})$$

282 with  $P^{99}$  denoting the percentile function for the 99<sup>th</sup> percentile. For precipitation, all-day percentiles  
283 (including the dry days) were used. Note the absolute nature of MAE99 and the fact that, in contrast  
284 to the BIAS metric, under- and overestimations of  $P^{99}$  at individual grid cells within a given sub-  
285 region do not compensate each other.

286 The similarity of the spatial pattern of climatological seasonal means was assessed using pattern  
287 correlation as defined by the Pearson product-moment coefficient of linear correlation

$$288 \quad \mathbf{PACO} = \frac{\text{cov}(\bar{X}_n, \bar{O}_n)}{\text{sd}(\bar{X}_n)\text{sd}(\bar{O}_n)}, \quad n = 1..N \quad (\text{Eq. 3})$$

289 with  $\text{cov}$  and  $\text{sd}$  representing the spatial covariance and standard deviation, respectively.

290 The interannual variability of seasonal means was evaluated using the ratio of interannual variability  
291 (RIAV). The spatial and temporal means of a season were first calculated for every year separately,  
292 and the standard deviations were then related according to

$$293 \quad \mathbf{RIAV} = \frac{\text{sd}(\bar{X}_y)}{\text{sd}(\bar{O}_y)} \quad (\text{Eq. 4})$$

294 The Cramér-von Mises Test (CMT; Anderson, 1962, Lunneborg, 2005) was used to evaluate the  
295 similarity of the cumulative distribution functions of daily values. In the case of precipitation, only the  
296 wet days were considered (wet-day threshold of 1mm day<sup>-1</sup>). In order to remove the influence of the  
297 bias in the mean (which is evaluated already by the BIAS metric) the climate model data were first  
298 corrected for the mean bias. For temperature and precipitation this was done by additive and  
299 multiplicative correction, respectively. After the bias correction, the CMT was applied to every grid  
300 point separately resulting in a probability value for rejection  $p_n$ . Using a significance level of 0.05, the  
301 fraction of grid-points with non-rejection (i.e., the null-hypothesis of the two distributions being  
302 similar cannot be rejected at a probability of 0.05) was calculated. The latter represents the final  
303 Cramér-von Mises performance metric  $CM$ . In mathematical terms, this can be described as follows:

$$304 \quad p_n = \text{CMT}(X_n, O_n) \quad (\text{Eq. 5})$$

$$305 \quad \begin{aligned} c_n &= 1 & \text{if } p_n > 0.05 \\ c_n &= 0 & \text{if } p_n \leq 0.05 \end{aligned} \quad (\text{Eq. 6})$$

$$306 \quad \mathbf{CM} = \frac{1}{N} \sum_{n=1}^N c_n \quad (\text{Eq. 7})$$

307 Note that this simple version of the metric neglects a potential spatial autocorrelation of the test  
308 statistic and does not consider field significance (e.g., Ivanov et al. 2017a and 2017b). Two further  
309 metrics were only calculated for either temperature or precipitation. For temperature only, the mean  
310 absolute error of the 1<sup>st</sup> percentile was used to evaluate moderately cold extremes:

$$311 \quad \mathbf{MAE01} = \frac{1}{N} \sum_{n=1}^N |P^1(X_n) - P^1(O_n)| \quad (\text{Eq. 8})$$

312 For precipitation only, the mean absolute bias in the wet-day frequency was evaluated by

$$313 \quad \mathbf{WDFREQ} = \frac{1}{N} \sum_{n=1}^N |wdf r(X_n) - wdf r(O_n)| \quad (\text{Eq. 9})$$

314 with  $wdf r()$  being the wet-day frequency [%] for a given grid point and a given season for a wet-day  
315 threshold of 1 mm day<sup>-1</sup>.

## 316 2.5 Uncertainty Intercomparison

317 A dedicated comparison framework was employed to quantify the relation between *observational*  
318 *uncertainty* (the influence of the choice of the reference dataset on the evaluation) and *model*  
319 *uncertainty* (the effect of the choice of a specific RCM on the evaluation). In case observational  
320 uncertainty is large, model evaluation against one specific reference dataset has to be considered as  
321 non-robust and evaluation exercises need to definitely take into account observational uncertainty.

322 Let  $P_{i,j}$  be the value of given performance metric for a given variable, sub-region and season when  
323 employing reference dataset  $i$  ( $i \in \{1,2,3\}$ ) for evaluating RCM  $j$  ( $j \in \{1,2,3,4,5\}$ ). Observational  
324 uncertainty is defined as the mean standard deviation of the metric's values when comparing an  
325 RCM against each of the three reference datasets:

$$326 \quad U_{OBS} = \frac{\sum_{j=1}^5 \sqrt{\frac{1}{2} \sum_{i=1}^3 (P_{i,j} - \frac{1}{3} \sum_{i=1}^3 P_{i,j})^2}}{5} \quad (\text{Eq. 10})$$

327 Correspondingly, model uncertainty is defined as the mean standard deviation of the respective  
328 metric's values when comparing all  $\binom{5}{3} = 10$  three-member RCM sub-ensembles against a given  
329 reference dataset:

$$330 \quad U_{MOD} = \frac{\sum_{i=1}^3 (\frac{1}{10} \sum_{n=1}^{10} \sqrt{\frac{1}{2} \sum_{j \in S_n} (P_{i,j} - \frac{1}{3} \sum_{j \in S_n} P_{i,j})^2})}{3} \quad (\text{Eq. 11})$$

331 where  $S_n = \{(1,2,3); (1,2,4); (1,2,5); (1,3,4); (1,3,5); (1,4,5); (2,3,4); (2,3,5); (2,4,5); (3,4,5)\}$ .  
332 Three-member sub-ensembles are chosen to be consistent with  $U_{obs}$ . The ratio

$$333 \quad R = \frac{U_{OBS}}{U_{MOD}} \quad (\text{Eq. 12})$$

334 for a given metric, variable, sub-region and season then defines the ratio of observational and model  
335 uncertainty. If this ratio is larger than 1 observational uncertainty is larger than model uncertainty  
336 and, hence, presents an important contribution to overall evaluation uncertainty and should be  
337 considered in evaluation exercises. Note that in our case model uncertainty is defined via the spread  
338 among different re-analysis driven RCMs. When evaluating RCM experiments that are driven by  
339 different GCMs at their lateral boundaries (i.e. the kind of experiments employed for regional climate  
340 projections) this spread and, hence, model uncertainty can be expected to be larger.

341 As mentioned earlier, the observational references except for HR over Sweden have not been  
342 corrected for precipitation undercatch and might underestimate true precipitation sums which can  
343 have an effect on the uncertainty intercomparison. In a dedicated sensitivity analysis we therefore  
344 carried out a modified uncertainty analysis for precipitation. For this purpose, a bulk correction of 20  
345 % was applied to all observational references (E-OBS, MESAN and HR, except for HR over sub-region  
346 SW), i.e. daily precipitation amounts were multiplied by a factor of 1.2. This bulk correction might  
347 underestimate the undercatch in winter in some regions and overestimate it in summer. It should  
348 only be considered as a rough estimate employed to address the principle sensitivity of our

349 uncertainty analysis with respect to the undercatch issue. Simulated precipitation amounts were not  
 350 modified. Uncertainty ratios  $R$  were re-computed employing the undercatch-corrected observations.

## 351 2.6 Ranking Framework

352 As model selection and weighting schemes (e.g., Christensen et al., 2010) are commonly based on  
 353 the assessment of a climate model's ability to simulate the present-day climate (Räisänen et al.,  
 354 2007), part of the uncertainty in these schemes arises from differences between the available  
 355 reference datasets. To test how the relative performance of the RCMs depends on the selected  
 356 reference, a simple scheme combining the performance metrics introduced in Section 2.4 was used.  
 357 First, to ensure that smaller values indicate better RCM performance absolute values were  
 358 considered for BIAS while RIAV, PACO and CM were transformed according to

$$359 P' = |1 - P| \quad (\text{Eq. 13})$$

360 with  $P$  being the value of the respective performance metric. MAE99, MAE01 and WDFREQ were  
 361 used as computed according to Eqs. 2, 8 and 9, respectively. For a consistent combination of the  
 362 metrics the values were furthermore normalized (Santer et al., 2009; Rupp et al., 2013) to obtain the  
 363 respective score  $S \in [0,1]$  for a given model  $j$  and performance metric  $m$  (indices for season  $k$ ,  
 364 region  $r$  and reference dataset  $i$  omitted):

$$365 S_{j,m} = 1 - \frac{P_{j,m} - \min(P_m)}{\max(P_m) - \min(P_m)} \quad (\text{Eq. 14})$$

366 with  $\min(P_m)$  and  $\max(P_m)$  denoting the minimum/maximum value of the five  $P_{j,m}$  (five RCMs  $j$ ) for  
 367 the case considered. Note that, in contrast to the performance metric  $P_{j,m}$ , the larger the value of  
 368 the score  $S_{j,m}$  the better the performance of a particular RCM for a given performance metric. For  
 369 each reference dataset and each variable, the final overall normalized scores were then calculated  
 370 separately for each RCM  $j$  and region  $r$  by taking an average over  $K$  seasons and  $M$  performance  
 371 metrics:

$$372 \bar{S}_{j,r} = \frac{1}{MK} \sum_{m=1}^M \sum_{k=1}^K S_{j,k,m,r} \quad (\text{Eq. 15})$$

373 Thus, equal weight is given to each performance score. The RCM simulations were then ranked  
 374 according to the obtained  $\bar{S}_{j,r}$  values separately for temperature and precipitation ( $M=6$  in Eq. 15). If  
 375 there are no systematic differences in the relative RCM performance (i.e.,  $S_{j,k,m,r}$  tends to vary  
 376 randomly for a given model  $j$  and a given region  $r$ )  $\bar{S}_{j,r}$  is expected to approach 0.5. Combined  
 377 temperature and precipitation ranks were computed by considering both temperature and  
 378 precipitation metrics in Eq. 15 ( $M=12$ ).

379 A similar scheme with a slightly different set of performance metrics was compared to a more  
 380 sophisticated scheme by Rupp et al. (2013) and was found to yield qualitatively similar results. One  
 381 should note that model ranking is inherently subjective (Overland et al., 2011) and depends on the  
 382 selected climatic aspects, error measures as well as the temporal and spatial scales considered.  
 383 However, for illustrational purposes the selected scheme is considered sufficient.

## 384 3. Reference Data Uncertainty

385 In order to provide a first impression on the differences among the three reference datasets which  
 386 will ultimately determine differences in the RCM evaluation exercise we here present a comparison  
 387 of E-OBS, MESAN and HR in terms of the spatial distribution of climatological seasonal mean values.  
 388 This comparison is directly relevant for the BIAS metric but might also concern metrics such as

389 MAE99, MAE01 or PACO. For the comparison we assume HR as reference (due to its highest  
390 underlying network density) and display the differences of E-OBS and MESAN with respect to HR.

391 Figure 2 shows the spatial distribution of seasonal mean temperature in HR and the corresponding  
392 deviations of E-OBS and MESAN. All three datasets agree on the general continental-scale  
393 temperature gradients and on large-scale mean values (not explicitly shown but deducible from  
394 Figure 2). Differences, however, appear over individual sub-regions and are obviously connected to  
395 the merging of different regional grids in the HR dataset and to complex orography. Over the Spanish  
396 Highlands, the Scandinavian Alps, Switzerland, south-western France and the Carpathians both E-OBS  
397 and MESAN can considerably deviate from HR in both seasons. Over the Carpathians these  
398 differences are systematic in the sense that HR provides the highest temperatures. When moving to  
399 Poland, i.e. into a region covered by a different sub-regional dataset in HR, a close agreement  
400 between E-OBS, MESAN and HR is obtained in both seasons. Over south-western France, in contrast,  
401 HR systematically shows lower temperatures. Over most parts of Spain MESAN yields lower winter  
402 temperatures than HR with differences partly larger than 2°C. Again, this bias pattern disappears  
403 when moving into France where mean winter temperatures in HR and MESAN closely agree. E-OBS  
404 shows the highest temperatures over FR in both seasons with differences to HR often larger than  
405 0.5°C. This might be connected to the fact that most French station data underlying E-OBS represents  
406 larger urban settings possibly affected by the urban heat island effect (see, e.g.,  
407 <http://www.ecad.eu/download/stations.txt>). Further consistent features are higher temperatures in  
408 E-OBS and MESAN over parts of Scandinavia and the European Alps.

409 Regarding mean seasonal precipitation all reference datasets again agree on the basic continental  
410 scale patterns and on large-scale mean values (not explicitly shown but deducible from Figure 3). A  
411 noticeable difference in comparison to HR is the general underestimation of precipitation by both E-  
412 OBS and MESAN in both seasons and for most parts of the analysis domain. Deviations can be as  
413 large as 50% (e.g. Poland and Sweden in MESAN with respect to HR). Exceptions are the complex  
414 coastline of Western Norway, where E-OBS provides higher precipitation sums than HR in both  
415 winter and summer, and Spain, where MESAN precipitation is comparable to HR and over parts of  
416 the country even higher in summer. The same is true over parts of the Carpathians. Over France,  
417 MESAN and HR are in very close agreement, which is likely connected to the good station coverage in  
418 MESAN over this region and which supports findings by Isotta et al. (2015) and Prein and Gobiet  
419 (2017) for (south-eastern) France. The general picture of highest precipitation sums in HR and drier  
420 conditions in E-OBS and MESAN might be a direct consequence of the higher underlying network  
421 density in HR and the fact that more high-elevation stations are sampled. Over Sweden, a further  
422 reason for lower precipitation sums in E-OBS and MESAN especially in wintertime is presumably the  
423 applied undercatch correction in the PTHBV dataset underlying HR in this region (see Section 2.1.2).

## 424 4. Model Evaluation

425 In the following the results of the model evaluation exercise are presented separately for both  
426 variables (temperature and precipitation) and for both seasons (winter and summer). The analysis  
427 allows for a separate assessment of each performance metric, each observational reference and each  
428 sub-region. For the sake of clarity and according to the objectives of this work we do not explicitly  
429 identify the five individual RCM experiments (see Table 1 for their identification though).

### 430 4.1 Temperature

431 Figures 4 and 5 present the temperature evaluation results for the winter and the summer season,  
432 respectively. In most cases the spatially averaged model biases (BIAS) approximately agree for the  
433 three reference datasets. In winter (Fig. 4) a cold model bias prevails and, depending on sub-region  
434 and RCM, can amount to more than -2 °C. The range of model biases (given by the vertical extent of

435 the bars) is largest over sub-region CH, where two of the five RCMs are subject to pronounced cold  
436 biases. This is likely related to the strong topographic variability of this domain, the pronounced  
437 differences in RCM orographies and to the fact that the applied lapse rate correction is based on the  
438 simplifying assumption of a global lapse rate being stationary in both time and space. Over France  
439 and Norway cold biases are most pronounced when evaluating against E-OBS, which is in line with  
440 the higher winter temperatures in E-OBS compared to HR and MESAN over these sub-domains (cf.  
441 Fig. 2).

442 In the summer season (Fig. 5) notable differences of the BIAS metric when comparing against  
443 different reference datasets are apparent for sub-regions CA, FR, NO, SP and CH. For the other  
444 regions, spatially averaged model biases mostly agree. A similar finding is obtained for MAE99 and  
445 MAE01. In the latter case, however, the evaluation against MESAN yields considerably larger summer  
446 model biases than for E-OBS and HR in the topographically complex sub-regions NO and CH. Note the  
447 extremely large ranges of MAE01 in the winter season with differences between the RCMs of more  
448 than 9 °C in sub-region CH, independently of the reference dataset. These large biases are found in  
449 two of the five RCMs only. In general, the large model spread over CH indicates difficulties of the  
450 RCMs to reliably reproduce minimum temperatures over regions of complex topography.

451 Concerning the RIAV metric the evaluation results are robust with respect to the choice of reference  
452 in both winter and summer with minor exceptions only. Model uncertainty as expressed by the  
453 vertical extent of the bars is generally much larger than the influence of the reference dataset. The  
454 situation is different though for the PACO metric. Here, the choice of the reference can have an  
455 important influence on the evaluation results. Correlation coefficients are high in general (> 0.8 in all  
456 cases) owing mainly to the pronounced influence of topography on spatial temperature patterns  
457 which is, in principle, represented by both the RCMs and by the observations. Depending on sub-  
458 region and season, reference data uncertainty can however strongly dominate. Use of the HR dataset  
459 as reference leads to lower correlation coefficients in winter in sub-regions FR and SP. The same is  
460 true for FR, PO and SP in summer. These results suggest differences in the spatial pattern of seasonal  
461 mean temperatures in the three reference datasets even for regions of pronounced topography and  
462 even for the aggregated evaluation scale of 0.22°.

463 For the distribution-based CM metric, the choice of the reference has an important effect in a few  
464 cases only and model uncertainty mostly dominates. The choice of the reference dataset markedly  
465 influences CM over CA, FR and SP in winter and CA and CH in summer.

## 466 4.2 Precipitation

467 For precipitation, a pronounced dependency of the BIAS metric on the choice of the reference can be  
468 found in both seasons but depending on the sub-region (Figs. 6 and 7). In winter and for sub-region  
469 SP, positive model biases with respect to E-OBS can partly translate into negative biases with respect  
470 to HR, reflecting the higher precipitation sums in HR compared to E-OBS over most parts of sub-  
471 region SP (cf. Figure 3). The same is true for SW and CH in summer. In a few cases the BIAS ranges for  
472 the three reference datasets only slightly overlap and reference data uncertainty is of a similar  
473 magnitude as model uncertainty (for instance, sub-regions CA, PO and SW in winter). In the last case  
474 (SW in winter) a possible reason is the undercatch correction of the Swedish HR dataset that  
475 potentially reduces positive model biases compared to the non-corrected E-OBS and MESAN data.

476 For MAE99, i.e. for the upper tail of the daily precipitation distribution, reference data uncertainty  
477 has a larger magnitude than for the BIAS metric (note the different y-axis scales in the upper left and  
478 upper middle panels) but is clearly dominated by model uncertainty, especially in summertime. A  
479 completely different result is obtained for the spatially averaged absolute bias of the wet day  
480 frequency WDFREQ. While model biases with respect to E-OBS and HR approximately agree, the use  
481 of the MESAN reanalysis as reference is in most cases associated with larger biases that are partly  
482 outside the bias range obtained for E-OBS and HR. The reason is a considerably lower wet day

483 frequency in MESAN compared to E-OBS and HR and a generally positive wet day frequency bias of  
484 the RCMs. This bias, and hence WDFREQ, is therefore largest when using MESAN as reference.

485 In wintertime and over sub-regions PO, SW and CH the MESAN reanalysis is furthermore associated  
486 with larger RIAV values (Fig. 6, lower left panel), i.e. a more pronounced overestimation of  
487 interannual precipitation variability. All other cases show similar RIAV ranges regardless of the  
488 reference employed and model uncertainty clearly dominates. For PACO the results considerably  
489 depend on the sub-region. As a general picture, PACO values are systematically lower compared to  
490 temperature which reflects the less pronounced control of topography on the spatial pattern of  
491 mean seasonal precipitation. The PACO ranges for the three reference datasets are similar in many  
492 cases but there are exceptions. The use of E-OBS, for instance, leads to considerably lower values  
493 over sub-regions CA, FR and SP in winter while HR is associated with a lower pattern correlation for  
494 PO but higher values for sub-region SW. In summer, MESAN is associated with lower correlations  
495 over CA, and HR with higher correlations over CA and SW. Overall, however, model uncertainty  
496 dominates for the PACO metric.

497 A different picture is obtained for the distribution-based CM metric (lower right panels). The range of  
498 CM values for a given reference dataset is generally high, but especially the use of MESAN as  
499 reference can be associated with much lower values compared to E-OBS and HR, i.e. with a lower  
500 fraction of grid cells passing the CM test. This feature affects all sub-regions in winter and sub-  
501 regions GE, NO, PO, SW and CH in summer. It is obviously associated with the much higher WDFREQ  
502 value when using MESAN as a reference, i.e. with the lower wet day frequency in MESAN. Note that  
503 CM only considers the wet-day distribution (see Section 2.4) and is not directly affected by wet-day  
504 frequency biases. The close relation between both metrics hence indicates that model biases in the  
505 wet-day frequency when comparing against MESAN come along with biases in the precipitation  
506 distribution for wet days only, i.e. that at least the complete lower tails of the two all-day  
507 distributions (model and reference) considerably differ from each other.

## 508 **5. Observational Versus Model Uncertainty**

509 We here present the results of the uncertainty intercomparison introduced in Section 2.5. This  
510 analysis can be seen as a summary of the comparison between observational uncertainty (offset of  
511 the three vertical bars for a given performance metric, season and sub-region) and model uncertainty  
512 (vertical extent of the bars) provided in Chapter 4 and apparent from Figures 4 to 7.

513 For temperature (Fig. 8) uncertainty ratios smaller than one are obtained in most cases, i.e.  
514 observational uncertainty is typically smaller than model uncertainty. But exceptions to this general  
515 pattern are possible, and also the magnitude of the uncertainty ratio primarily depends on the  
516 performance metric considered. For the seasonal mean model bias (BIAS) ratios are consistently  
517 smaller than 0.5, indicating a model uncertainty being twice as large as observational uncertainty.  
518 With reference to the scores describing the tails of the daily values (MAE99 and MAE01) and the  
519 frequency distribution (CM), observational uncertainty is also smaller than model uncertainty with  
520 the exception of Spain for MAE01 in summer. Ratios for RIAV are below one throughout all sub-  
521 regions and both seasons with typically somewhat larger values in winter. In contrast to all other  
522 performance metrics, the ratios for the pattern correlation PACO are close to or larger than one in at  
523 least half of the cases, i.e. observational uncertainty dominates. This is in particular true for sub-  
524 regions FR, GE and SP during summer.

525 As a general pattern observational uncertainty, i.e. the choice of the reference data, tends to be  
526 more important for precipitation (Fig. 9) than for temperature. Uncertainty ratios for WDFREQ and  
527 CM are close to or even larger than one in most cases. Maximum values larger than three are  
528 obtained for winter in sub-regions CA, PO and SW (WDFREQ) and for sub-region CA (CM). For the  
529 cases of WDFREQ and CM these high values are clearly related to low WDFREQ values in the MESAN

530 reference, which constitute an outlier within the observational ensemble. They are probably related  
531 to specifics of the MESAN spatial interpolation and not to shortcomings in the underlying station  
532 observations. Except for a few cases summer ratios are smaller than their winter counterparts,  
533 indicating a smaller contribution of observational uncertainty in summer. This is mainly due to the  
534 fact that MESAN deviates stronger from E-OBS and HR in winter than in summer. A clearly  
535 dominating observational uncertainty is also found for PACO in sub-region CA (both seasons) as well  
536 as in PO and SW (winter only). The same is true for the winter BIAS in sub-regions CA, PO and SW.  
537 The latter are, however, outliers since for the BIAS metric ratios close to 0.5 are obtained for most  
538 other cases, i.e. model uncertainty clearly dominates. Also for MAE99 and RIAV ratios smaller than  
539 one are obtained with the exception of CA (MAE99) and PO and SW (RIAV) in winter

540 The influence of a potential precipitation undercatch in the observational references on the  
541 uncertainty analysis can be derived from Figure 10. For most performance metrics the uncertainty ratios  
542 are not or only slightly modified compared to the original results (Fig. 9). The most important change  
543 is obtained for the MAE99 metric which is especially sensitive as it considers absolute biases at the  
544 upper tail of the distribution. Here, uncertainty ratios are increasing in many cases. Roughly, the  
545 same stands for further measures based on daily data (WDFREQ and CM). For sub-region SW,  
546 specifically, observational uncertainty for MAE99 grows in both winter and summer as only two of  
547 the references (E-OBS and MESAN) were corrected for the undercatch compared to the original  
548 analysis. This results in larger inter-observational differences, in a larger observational uncertainty  
549 and, hence, in a larger uncertainty ratio. For the BIAS metric over SW in winter, undercatch  
550 correction brings the three references closer together (not shown), resulting in a decreasing  
551 observational uncertainty and a decreasing uncertainty ratio.

## 552 6. Model Ranking

553 To assess the influence of observational uncertainty on model ranking we first show the results for  
554  $S_{j,m}$  (Eq. 14; simply denoted as  $S$  hereafter) separately for temperature and precipitation when  
555 averaged over all seasons and regions (Fig. 11). For illustrational purposes, the actual RCM ranks  
556 based on  $S$  are also shown. The individual performance metrics show a varying degree of variation in  
557  $S$  between the reference datasets (horizontal variation within a given panel). BIAS and MAE99 have a  
558 similar normalized error pattern and almost identical ranks for all reference datasets. Model C, for  
559 instance has the best performance for both metrics in terms of temperature, independently of the  
560 reference dataset. In contrast, model D shows the worst performance for temperature but the best  
561 for precipitation. In contrast to these cases of agreement between reference datasets, scores for CM  
562 (precipitation) and PACO (temperature) show noticeable differences when employing E-OBS, MESAN  
563 or HR as reference. Unsurprisingly, variations are even larger when individual regions are considered  
564 (not shown). Concerning the performance of a given model for different performance metrics  
565 (vertical variation within a given panel) model C, for instance, has the highest  $S$  values (and the best  
566 ranking) for most temperature performance metrics, while model D shows the best performance in  
567 the case of precipitation. While not the worst performing model in all cases, model E often shows the  
568 lowest  $S$  and ranks poorly accordingly, regardless of the reference dataset considered. The fact that  
569 the dependence of the evaluation results on the reference dataset in turn depends on the metric  
570 considered confirms findings from previous works (e.g., Santer et al., 2009, Rupp et al., 2013) and  
571 should be kept in mind when interpreting the ranking results.

572 To illustrate the results for the full ranking scheme, Fig. 12 presents the overall normalized score  $\bar{S}_{j,r}$   
573 of Eq. 15 (denoted as  $\bar{S}$  hereafter) for each sub-region together with the actual RCM ranks. As an  
574 overall picture RCM ranks are similar, independently of the reference dataset employed. However,  
575 differences in  $\bar{S}$  between the reference datasets can be non-negligible depending on the region and  
576 RCM considered. On average, differences in  $\bar{S}$  between the reference datasets are largest over sub-  
577 regions GE and PO, although individual RCMs also stand out in other regions such as SP (model C) or



578 FR (model D). On the other hand, Switzerland (CH) shows only small differences in the overall scores  
579 between the reference datasets. Furthermore, variations in the actual ranks depending on the  
580 reference dataset employed are apparent. These differences tend to be smallest in CH, NO and SW,  
581 where the intermodel differences in  $\bar{S}$  are relatively large compared to the differences between the  
582 reference datasets. In other sub-regions a change of the reference dataset can lead to larger changes  
583 in the model ranks (e.g., the rank of model C in SP can change by four levels). This shows that model  
584 ranking becomes more dependent on the reference dataset when spatial details are considered.  
585 Finally, although the best performing RCM depends on the region and the reference, a noticeable  
586 feature is the systematically poor performance of model E in comparison to other models. Model E  
587 has the lowest rank in almost all cases regardless of the reference, and values of  $\bar{S}$  rarely approach  
588 0.5 for this model.

## 589 7. Summary and Conclusions

590 The objective of the present work was to illustrate the effect of uncertainties in gridded  
591 observational reference datasets on RCM evaluation for two variables (temperature and  
592 precipitation) on a pan-European scale. For this purpose we made use of three different gridded  
593 observational reference datasets (E-OBS v15, national/regional high-resolution grids (HR), EURO4M  
594 MESAN) and five reanalysis-driven RCM experiments carried out within the EURO-CORDEX initiative.  
595 Our well-defined performance assessment framework considers a range of performance metrics for  
596 eight different sub-regions of the European continent and includes an illustrative model ranking  
597 scheme. Note that the ensemble of reference grids is an ensemble of opportunity and is likely subject  
598 to inter-dependencies arising, for instance, from the use of common station time series in the  
599 interpolation or assimilation procedure. In general, an extension of the observational ensemble by,  
600 for instance, satellite-based products or by new upcoming datasets could alter the derived  
601 observational uncertainties and, hence, the overall evaluation results. The same would be true for an  
602 extension of the set of RCMs considered or for a different sampling of available RCM experiments.

603 A comparison of climatological seasonal mean values as represented by the three reference grids  
604 alone yields a general agreement concerning the continental-scale patterns, but also differences on  
605 regional scales. These depend on the variable, region and season considered and translate into  
606 differences in RCM performance scores. Largest differences in seasonal mean temperature occur  
607 over regions of pronounced topography, such as Spain, the European Alps, Scandinavia and the  
608 Carpathians. Except for the latter case, the high-resolution HR dataset typically shows lowest  
609 temperatures which might be related to a better sampling of high-elevation stations by HR. For the  
610 case of precipitation both MESAN and E-OBS typically underestimate mean seasonal precipitation as  
611 provided by HR.

612 For most performance metrics and especially for temperature, the influence of the choice of the  
613 observational reference on model evaluation is rather weak and is smaller than model uncertainty.  
614 This is especially true for winter temperature, where only the pattern correlation (PACO) and to  
615 some extent the distribution-based Cramér-von Mises score (CM) show notable dependencies on the  
616 reference dataset employed. However, winter PACO values are still larger than 0.8 for each individual  
617 sub-region and for any combination of RCM and reference dataset. Hence, spatial temperature  
618 patterns are, in a general sense, well represented by the RCMs independently of the specific  
619 reference employed. The same is true for the summer season which, however, is subject to slightly  
620 larger reference data uncertainty for PACO and for the mean absolute error of the 1<sup>st</sup> daily  
621 percentile (MAE01). For precipitation the influence of observational uncertainty is larger than for  
622 temperature. It often dominates model uncertainty especially for the absolute bias in the wet-day  
623 frequency (WDFREQ) and for the Cramér-von Mises score (CM) in winter. But even the spatially  
624 averaged measures of seasonal mean bias (BIAS) and ratio of interannual variability (RIAV) can be  
625 considerably affected by the choice of the reference observational product. The fact that most

626 observational references are not corrected for rain gauge undercatch has some influence on the final  
627 uncertainty analysis but does not change the general picture. Note that observational uncertainty  
628 being smaller than model uncertainty does not necessarily imply that uncertainties in observations  
629 are negligible and without influence. They can still be relevant, for instance, in model development  
630 or model bias correction.

631 When employing a simple and illustrative model ranking scheme on these results it is found that RCM  
632 ranking can depend on the reference dataset employed, and more often for precipitation than for  
633 temperature. In individual cases, final model ranks can differ by up to four (out of five models)  
634 depending on the choice of the reference dataset. These findings are in line with previous works  
635 (e.g., Gómez-Navarro et al., 2012; Prein and Gobiet, 2016) which suggests that uncertainties related  
636 to the reference data should ideally be taken into account when assessing climate model  
637 performance in the present-day climate. However, if a focus is laid on temperature only the three  
638 reference datasets agree to a large extent, indicating the suitability of each individual product for  
639 climate model evaluation purposes. Furthermore, spatio-temporally averaged temperature and  
640 precipitation climates are very similar among the three references (see the BIAS metric), and model  
641 uncertainty clearly dominates in these case. Also note that all datasets employed in the present work  
642 were aggregated to the comparatively low E-OBS grid spacing of 0.22° prior to the analysis, including  
643 the high-resolution HR data. This spatial aggregation might to some extent mask the added value of  
644 HR but is required in the context of the present work. The full benefits of the higher-resolved HR data  
645 and their underlying dense station network will however only become apparent when evaluating, for  
646 instance, very high resolution RCM experiments at the convection-resolving scale (e.g. Ban et al.,  
647 2014).

648 Considering the ranking exercise itself, one should keep in mind that the ranking scheme applied  
649 here is likely to suffer from commonly known limitations (Overland et al., 2011; Santer et al., 2009,  
650 Rupp et al., 2013) and that the results are specific for the selected RCMs and performance metrics.  
651 On the other hand, it has been previously shown that only small uncertainties in the ranking and  
652 weighting of models can result in strong differences and potentially misleading signals (Weigel et al.,  
653 2010).

654

#### 655 **Acknowledgments**

656 The present work has been carried out as part of the EU-COST Action VALUE (Validating and  
657 Integrating Downscaling Methods for Climate Change Research; ES1102). We gratefully acknowledge  
658 the providers of RCM and observational data. For the high-resolution national/regional grids these  
659 are the University of Cantabria (SP), the Institute of Meteorology and Water Management - National  
660 Research Institute (PO), Météo-France/CERFACS (FR), The Swedish Meteorological and Hydrological  
661 Institute (SE), Deutscher Wetterdienst (GE), the Hungarian Meteorological Service (CA), the  
662 Norwegian Meteorological Institute (NO) and Federal Office of Meteorology and Climatology  
663 MeteoSwiss (CH). Furthermore, we acknowledge the E-OBS dataset from the EU-FP6 project  
664 ENSEMBLES (<http://ensembles-eu.metoffice.com>) and the data providers in the ECA&D project  
665 (<http://eca.knmi.nl>). The MESAN dataset was provided by the Swedish Meteorological and  
666 Hydrological Institute. All analysis were performed on the computing infrastructure of the Swiss  
667 National Supercomputing Centre CSCS. We furthermore thank the climate modelling groups of the  
668 EURO-CORDEX initiative for producing and making available their model output. The contribution of  
669 Olle Räty was partly funded by the Vilho, Yrjö and Kalle Väisälä Foundation of the Finnish Academy of  
670 Science and Letters.

671 **References**

- 672 Adam JC, Lettenmaier DP. 2003. Adjustment of global gridded precipitation for systematic bias.  
673 *Journal of Geophysical Research* 108: 4257 (D9). doi: 10.1029/2002JD002499.
- 674 Addor N, Fischer EM. 2015. The influence of natural variability and interpolation errors on bias  
675 characterization in RCM simulations. *Journal of Geophysical research – Atmospheres*, 120. doi:  
676 10.1002/2014JD022824.
- 677 Anderson TW. 1962. On the Distribution of the Two-Sample Cramer-von Mises Criterion. *The Annals*  
678 *of Mathematical Statistics* 33(3): 1148-1159. doi:10.1214/aoms/1177704477.
- 679 Awange JL, Ferreira VG, Forootan E, Khandu, Andam-Akorful SA, Agutu NO, He XF. 2016.  
680 Uncertainties in remotely sensed precipitation data over Africa. *International Journal of Climatology*  
681 36: 303-323. doi: 10.1002/joc.4346.
- 682 Ban N, Schmidli J, Schär C. 2014. Evaluation of the convection-resolving regional climate modeling  
683 approach in decade-long simulations. *Journal of Geophysical Research* 119 (13): 7889-7907. doi:  
684 10.1002/2014JD021478.
- 685 Begert M, Schlegel T, Kirchhofer W. 2005. Homogeneous temperature and precipitation series of  
686 Switzerland from 1864 to 2000. *International Journal of Climatology* 25: 65-80. doi:  
687 10.1002/joc.1118.
- 688 Beguería S, Vicente-Serrano SM, Tomás-Burguera M, Maneta M. 2016. Bias in the variance of gridded  
689 data sets leads to misleading conclusions about changes in climate variability. *International Journal of*  
690 *Climatology* 36(9): 3413-3422. doi: 10.1002/joc.4561.
- 691 Bellprat O, Kotlarski S, Lüthi D, Schär C. 2012. Exploring Perturbed Physics Ensembles in a Regional  
692 Climate Model. *Journal of Climate* 25: 4582-4599. doi: 10.1175/JCLI-D-11-00275.1.
- 693 Berg P, Norin L, Olsson J. 2016. Creation of a high resolution precipitation data set by merging  
694 gridded gauge data and radar observations for Sweden. *Journal of Hydrology* 541: 6-13. doi:  
695 10.1016/j.jhydrol.2015.11.031.
- 696 Brienen S, Früh B, Walter A, Trusilova K, Becker P. 2016. A Central European precipitation climatology  
697 – Part II: Application of the high-resolution HYRAS data for COSMO-CLM evaluation. *Meteorologische*  
698 *Zeitschrift* 25(2): 195-214. doi: 10.1127/metz/2016/0617.
- 699 Bucchignani E, Montesarchio M, Zollo AL, Mercogliano P. 2016. High-resolution climate simulations  
700 with COSMO-CLM over Italy: performance evaluation and climate projections for the 21st century.  
701 *International Journal of Climatology* 36: 735-756. doi: 10.1002/joc.4379.
- 702 Casanueva A, Herrera S, Fernández J, Frías MD, Gutiérrez JM. 2013. Evaluation and projection of daily  
703 temperature percentiles from statistical and dynamical downscaling methods. *Natural Hazards and*  
704 *Earth System Sciences* 13: 2089-2099. doi: 10.5194/nhess-13-2089-2013.
- 705 Cheneka BR, Brienen S, Fröhlich K, Asharaf S, Früh B. 2016. Searching for an added value in  
706 downscaled seasonal hindcasts over East Africa: COSMO-CLM forced by MPI-ESM. *Advances in*  
707 *Meteorology* 2016: 1-17. doi: 10.1155/2016/4348285.
- 708 Cheval S, Baciú M, Dumitrescu A, Breza T, Legatesb DR, and Chende V. 2010. Climatologic  
709 adjustments to monthly precipitation in Romania. *International Journal of Climatology* 31: 704-714.  
710 doi:10.1002/joc.2099.
- 711 Christensen JH, Kjellström E, Giorgi F, Lenderink G, Rummukainen M. 2010. Weight assignment in  
712 regional climate models. *Climate Research* 44: 179-194. doi: 10.3354/cr00916.

- 713 Dahlgren P, Gustafsson N. 2012. Assimilating host model information into a limited area model.  
714 *Tellus* 64A: 15836. doi: 10.3402/tellusa.v64i0.15836.
- 715 Davin EL, Maisonnave E, Seneviratne SI. 2016. Is land surface processes representation a possible  
716 weak link in current Regional Climate Models? *Environmental Research Letters* 11: 074027. doi:  
717 10.1088/1748-9326/11/7/074027.
- 718 Dee DP, Uppala SM, Simmons AJ, Berrisford P, Poli P, Kobayashi S, Andrae U, Balmaseda MA,  
719 Balsamo G, Bauer P, Bechtold P, Beljaars ACM, van de Berg L, Bidlot J, Bormann N, Delsol C, Dragani  
720 R, Fuentes M, Geer AJ, Haimberger L, Healy SB, Hersbach H, Hólm EV, Isaksen L, Kallberg P, Köhler M,  
721 Matricardi M, McNally AP, Monge-Sanz BM, Morcrette J-J, Park B-K, Peubey C, de Rosnay P, Tavolato  
722 C, Thépaut J-N, Vitart F. 2011. The ERA-Interim reanalysis: configuration and performance of the data  
723 assimilation system, *Quarterly Journal of the Royal Meteorological Society* 137: 553–597. doi:  
724 10.1002/qj.828.
- 725 Di Luca A, de Elía R, Laprise R. 2012. Potential for added value in precipitation simulated by high-  
726 resolution nested Regional Climate Models and observations. *Climate Dynamics* 38: 1229-1247. doi:  
727 10.1007/s00382-011-1068-3.
- 728 Dunn RJH, Donat MG, Alexander LV. 2014. Investigating uncertainties in global gridded data sets of  
729 climate extremes. *Climate of the Past* 10: 2171-2199. doi: 10.5194/cp-10-2171-2014.
- 730 Durand Y, Brun E, Mérindol L, Guyomarc'h G, Lesaffre B, Martin E. 1993. A meteorological estimation  
731 of relevant parameters for snow models. *Annals of Glaciology* 18: 65–71.
- 732 Frei C, Christensen JH, Déqué M, Jacob D, Jones RG, Vidale PL. 2003. Daily precipitation statistics in  
733 regional climate models: Evaluation and intercomparison for the European Alps. *Journal of*  
734 *Geophysical Research* 108(D3): 4124. doi: 10.1029/2002JD002287.
- 735 Frei C. 2014. Interpolation of temperature in a mountainous region using nonlinear profiles and non-  
736 Euclidean distances. *International Journal of Climatology* 34: 1585-1605. doi: 10.1002/joc.3786.
- 737 Frick C, Steiner H, Mazurkiewicz A, Riediger U, Rauthe M, Reich T, Gratzki A. 2014. Central European  
738 high-resolution gridded daily data sets (HYRAS): Mean temperature and relative humidity.  
739 *Meteorologische Zeitschrift* 23(1): 15-32. doi: 10.1127/0941-2948/2014/0560.
- 740 Fowler HJ, Blekinsop S, Tebaldi C. 2007. Linking climate change modelling to impacts studies: recent  
741 advances in downscaling techniques for hydrological modelling. *International Journal of Climatology*  
742 27: 1547-1578. doi: 10.1002/joc.1556.
- 743 Gandin LS. 1965. Objective analysis of meteorological fields (Ob"ektivnyi analiz meteorologicheskikh  
744 polei), Translated from Russian by the Israel Program for Scientific Translations, Jerusalem. *Quarterly*  
745 *Journal of the Royal Meteorological Society*. doi: 10.1002/qj.49709239320.
- 746 Gbambie ASB, Poulin A, Boucher MA, Arsenault R. 2017. Added value of alternative information in  
747 interpolated precipitation datasets for hydrology. *Journal of Hydrometeorology* 18: 247-264. doi:  
748 10.1175/JHM-D-16-0032.1.
- 749 Gervais M, Tremblay LB, Gyakum JR, Atallah E. 2014. Representing extremes in a daily gridded  
750 precipitation analysis over the United States: Impacts of station density, resolution, and gridding  
751 methods. *Journal of Climate* 27: 5201-5218. doi: 10.1175/JCLI-D-13-00319.1.
- 752 Gómez-Navarro JJ, Montávez JP, Jerez S, Jiménez-Guerrero P, Zorita E. 2012. What is the role of the  
753 observational dataset in the evaluation and scoring of climate models? *Geophysical Research Letters*  
754 39: L24701. doi: 10.1029/2012GL054206.

- 755 Groisman PY, Legates DR. 1994. The accuracy of United States precipitation data, *Bulletin of the*  
756 *American Meteorological Society* 75: 215–227.
- 757 Häggmark L, Ivarsson K-I, Gollvik S, Olofsson P-O. 2000. MESAN, an operational mesoscale analysis  
758 system. *Tellus* 52A: 2-20.
- 759 Hanel M, Buishand A. 2011. Analysis of precipitation extremes in an ensemble of transient regional  
760 climate model simulations for the Rhine basin. *Climate Dynamics* 36: 1135–1153.
- 761 Hartmann DL, Klein Tank AMG, Rusticucci M, Alexander LV, Brönnimann S, Charabi Y, Dentener FJ,  
762 Dlugokencky EJ, Easterling DR, Kaplan A, Soden BJ, Thorne PW, Wild M, Zhai PM. 2013. Observations:  
763 Atmosphere and Surface. In: *Climate Change 2013: The Physical Science Basis. Contribution of*  
764 *Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*  
765 *[Stocker TF, Qin D, Plattner G-K, Tignor M, Allen SK, Boschung J, Nauels A, Xia Y, Bex V, Midgley PM*  
766 *(eds.)]. Cambridge University Press: Cambridge, UK and New York, NY, USA.*
- 767 Haslinger K, Anders I, Hofstätter M. 2013. Regional climate modelling over complex terrain: an  
768 evaluation study of COSMO-CLM hindcast model runs for the Greater Alpine Region. *Climate*  
769 *Dynamics* 40: 511-529. doi: 10.1007/s00382-012-1452-7.
- 770 Haylock MR, Hofstra N, Klein Tank AMG, Klok EJ, Jones PD, New M. 2008. A European daily high-  
771 resolution gridded data set of surface temperature and precipitation for 1950–2006. *Journal of*  
772 *Geophysical Research* 113: D20119. doi:10.1029/2008JD010201.
- 773 Hegerl GC, Jones PD, Barnett PD. 2001. Effect of observational sampling error on the detection and  
774 attribution of anthropogenic climate change. *Journal of Climate*, 14: 198–207.
- 775 Herold N, Alexander LV, Donat MG, Contractor S, Becker A. 2016. How much does it rain over land?  
776 *Geophysical Research Letters* 43: 341-348. doi: 10.1002/2015GL066615.
- 777 Herrera S, Gutiérrez JM, Ancell R, Pons MR, Frías MD, Fernández J. 2012. Development and analysis  
778 of a 50-year high-resolution daily gridded precipitation dataset over Spain (Spain02). *International*  
779 *Journal of Climatology* 32: 74–85. doi:10.1002/joc.2256.
- 780 Herrera S, Fernández J, Gutiérrez JM. 2016. Update of the Spain02 gridded observational dataset for  
781 EURO-CORDEX evaluation: assessing the effect of the interpolation methodology. *International*  
782 *Journal of Climatology* 36: 900-908. doi:10.1002/joc.4391.
- 783 Herrera S, Jaczewski A, Kotlarski S, Gutiérrez JM, Soares PMM. 2017. Sensitivity analysis of  
784 observational gridded datasets to the density of stations and the interpolation methodology. In  
785 preparation.
- 786 Hofstra N, Haylock M, New M, Jones PD. 2009. Testing E-OBS European high-resolution gridded data  
787 set of daily precipitation and surface temperature. *Journal of Geophysical Research* 114: D21101. doi:  
788 10.1029/2009JD011799.
- 789 Hofstra N, New M, McSweeney C. 2010. The influence of interpolation and station network density  
790 on the distribution and extreme trends of climate variables in gridded data. *Climate Dynamics* 35:  
791 841–858.
- 792 Isotta FA, Vogel R, Frei C. 2015. Evaluation of European regional reanalyses and downscalings for  
793 precipitation in the Alpine region. *Meteorologische Zeitschrift* 24: 15-37. doi:  
794 10.1127/metz/2014/0584.

- 795 Ivanov M, Warrach-Sagi K, Wulfmeyer V. 2017a. Field significance of performance measures in the  
796 context of regional climate model evaluation. Part 1: temperature. Theoretical and Applied  
797 Climatology, in press. doi: 10.1007/s00704-017-2100-2.
- 798 Ivanov M, Warrach-Sagi K, Wulfmeyer V. 2017b. Field significance of performance measures in the  
799 context of regional climate model evaluation. Part 2: precipitation. Theoretical and Applied  
800 Climatology, in press. doi: 0.1007/s00704-017-2077-x.
- 801 Jacob D, Petersen J, Eggert B, Alias A, Christensen OB, Bouwer LM, Braun A, Colette A, Déqué M,  
802 Georgievski G, Georgopoulou E, Gobiet A, Menut L, Nikulin G, Haensler A, Hempelmann N, Jones C,  
803 Keuler K, Kovats S, Kröner N, Kotlarski S, Kriegsman A, Martin E, van Meijgaard E, Moseley C, Pfeifer  
804 S, Preuschmann S, Radermacher C, Radtke K, Rechid D, Rounsevell M, Samuelsson P, Somot S,  
805 Soussana J-F, Teichmann C, Valentini R, Vautard R, Weber B, Yiou P. 2014. EURO-CORDEX: new high-  
806 resolution climate change projections for European impact research. Regional Environmental Change  
807 14: 563-578. doi: 10.1007/s10113-013-0499-2.
- 808 Johansson B, Chen D. 2003. The influence of wind and topography on precipitation distribution in  
809 Sweden: statistical analysis and modelling. International Journal of Climatology 23(12): 1523–1535.
- 810 Kotlarski S, Block A, Böhm U, Jacob D, Keuler K, Knoche R, Rechid D, Walter A. 2005. Regional climate  
811 model simulations as input for hydrological applications: evaluation of uncertainties. Advances in  
812 Geosciences 5: 119-125.
- 813 Kotlarski S, Hagemann S, Krahe P, Podzun R, Jacob D. 2012. The Elbe river flooding 2002 as seen by  
814 an extended regional climate model. Journal of Hydrology 472-473: 169-183. doi:  
815 10.1016/j.jhydrol.2012.09.020.
- 816 Kotlarski S, Keuler K, Christensen OB, Colette A, Déqué M, Gobiet A, Goergen K, Jacob D, Lüthi D, van  
817 Meijgaard E, Nikulin G, Schär C, Teichmann C, Vautard R, Warrach-Sagi K, Wulfmeyer V. 2014.  
818 Regional climate modeling on European scales: a joint standard evaluation of the EURO-CORDEX RCM  
819 ensemble. Geoscientific Model Development 7: 1297-1333. doi: 10.5194/gmd-7-1297-2014.
- 820 Kyselý J, Plavcová E. 2010. A critical remark on the applicability of E-OBS European gridded  
821 temperature data set for validating control climate simulations. Journal of Geophysical Research 115:  
822 D23118. doi: 10.1029/2010JD014123.
- 823 Lakatos M, Szentimrey T, Bihari Z, Szalai S. 2013. Creation of a homogenized climate database for the  
824 Carpathian region by applying the MASH procedure and the preliminary analysis of the data. Idojaras  
825 117: 143–158.
- 826 Lenderink G. 2010. Exploring metrics of extreme daily precipitation in a large ensemble of regional  
827 climate model simulations. Climate Research 44: 151-166. doi: 10.3354/cr00946.
- 828 Lunneborg CE. 2005. Cramer-Von Mises Test. Wiley StatsRef: Statistics Reference Online. doi:  
829 10.1002/9781118445112.stat06554.
- 830 Lussana C, Tveito OE, Uboldi F. 2016. seNorge v2.0, Temperature: An observational gridded dataset  
831 of temperature for Norway. MET Report No. 14/2016. Available from  
832 [https://www.met.no/sokeresultat/\\_attachment/inline/243074f4-09bf-4f63-b98a-  
833 f329b3661ce4:c586d2b116d185dc2ac000a1eca6cd98f2f5bdbd/MET-report-14-2016.pdf](https://www.met.no/sokeresultat/_attachment/inline/243074f4-09bf-4f63-b98a-f329b3661ce4:c586d2b116d185dc2ac000a1eca6cd98f2f5bdbd/MET-report-14-2016.pdf) (last access:  
834 18 July 2017).
- 835 Mahlstein I, Spirig C, Liniger MA, Appenzeller C. 2015. Estimating daily climatologies for climate  
836 indices derived from climate model data and observations. Journal of Geophysical Research –  
837 Atmospheres 120: 2808-2818. doi: 10.1002/2014JD022327.

- 838 Maraun D, Osborn TJ, Rust HW. 2012. The influence of synoptic airflow on UK daily precipitation  
839 extremes. Part II: regional climate model and E-OBS data validation. *Climate Dynamics* 39: 287-301.  
840 doi: 10.1007/s00382-011-1176-0.
- 841 Maraun D, Wigmann M, Gutiérrez JM, Kotlarski S, Chandler RE, Hertig E, Wibig J, Huth R, Wilcke RAL.  
842 2015. VALUE: A framework to validate downscaling approaches for climate change studies. *Earth's*  
843 *Future* 3. doi: 10.1002/2014EF000259.
- 844 McMillan H, Krueger T, Freer J. 2012. Benchmarking observational uncertainties for hydrology:  
845 rainfall, river discharge and water quality. *Hydrological Processes* 26: 4078-4111. doi:  
846 10.1002/hyp.9384.
- 847 MeteoSwiss. 2013a: Daily Mean, Minimum and Maximum Temperature: TabsD, TminD ,TmaxD.  
848 Available from [www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-](http://www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-publikationen/produkt/raeumliche-daten-temperatur/doc/ProdDoc_TabsD.pdf)  
849 [publikationen/produkt/raeumliche-daten-temperatur/doc/ProdDoc\\_TabsD.pdf](http://www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-publikationen/produkt/raeumliche-daten-temperatur/doc/ProdDoc_TabsD.pdf) (last access: 18 July  
850 2017).
- 851 MeteoSwiss. 2013b: Daily Precipitation (final analysis): RhiresD. Available from  
852 [www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-](http://www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-publikationen/produkt/raeumliche-daten-niederschlag/doc/ProdDoc_RhiresD.pdf)  
853 [publikationen/produkt/raeumliche-daten-niederschlag/doc/ProdDoc\\_RhiresD.pdf](http://www.meteoswiss.admin.ch/content/dam/meteoswiss/de/service-und-publikationen/produkt/raeumliche-daten-niederschlag/doc/ProdDoc_RhiresD.pdf) (last access: 18  
854 July 2017).
- 855 Osborn T, Hulme M. 1997. Development of a relationship between station and grid-box rainyday  
856 frequencies for climate model evaluation. *Journal of Climate* 10: 1885-1908. doi: 10.1175/1520-0442.
- 857 Overland JE, Wang M, Bond NA, Walsh JE, Kattsov VM, Chapman WL. 2011. Considerations in the  
858 selection of global climate models for regional climate projections: the Arctic as a case study. *Journal*  
859 *of Climate* 24: 1583–1597. doi: 10.1175/2010JCLI3462.1.
- 860 Palazzi E, von Hardenberg J, Provenzale A. 2013. Precipitation in the Hindu-Kush Karakoram  
861 Himalaya: observations and future scenarios. *Journal of Geophysical Research: Atmospheres* 118: 85-  
862 100. doi: 10.1029/2012JD018697.
- 863 Prein AF, Gobiet A. 2017. Impacts of uncertainties in European gridded precipitation observations on  
864 regional climate analysis. *International Journal of Climatology* 37: 305-327. doi: 10.1002/joc.4706.
- 865 Quintana-Seguí P, Le Moigne P, Durand Y, Martin E, Habets F, Baillon M. 2008. Analysis of near-  
866 surface atmospheric variables: validation of the SAFRAN analysis over France. *Journal of Applied*  
867 *Meteorology and Climatology* 47: 92-107. doi: 10.1175/2007JAMC1636.1.
- 868 Rauthe M, Steiner H, Riediger U, Mazurkiewicz A, Gratzki, A. 2013. A Central European precipitation  
869 climatology? Part I: Generation and validation of a high-resolution gridded daily data set (HYRAS).  
870 *Meteorologische Zeitschrift* 22 (3): 235-256. doi: 10.1127/0941-2948/2013/0436.
- 871 Richter D. 1995. Ergebnisse methodischer Untersuchungen zur Korrektur des systematischen  
872 Messfehlers des Hellmann-Niederschlagsmessers. *Berichte des Deutschen Wetterdienstes* 194.  
873 Selbstverlag des Deutschen Wetterdienstes, Offenbach am Main. Available from  
874 [http://www.dwd.de/DE/leistungen/pbfb\\_verlag\\_berichte/pdf\\_einzelbaende/194\\_pdf.pdf](http://www.dwd.de/DE/leistungen/pbfb_verlag_berichte/pdf_einzelbaende/194_pdf.pdf) (last  
875 access: 18 July 2017).
- 876 Ring C, Mannig B, Pollinger F, Paeth H. 2016. Uncertainties in the simulation of precipitation in  
877 selected regions of humid and dry climate. *International Journal of Climatology* 36: 3521-3538. doi:  
878 10.1002/joc.4573.

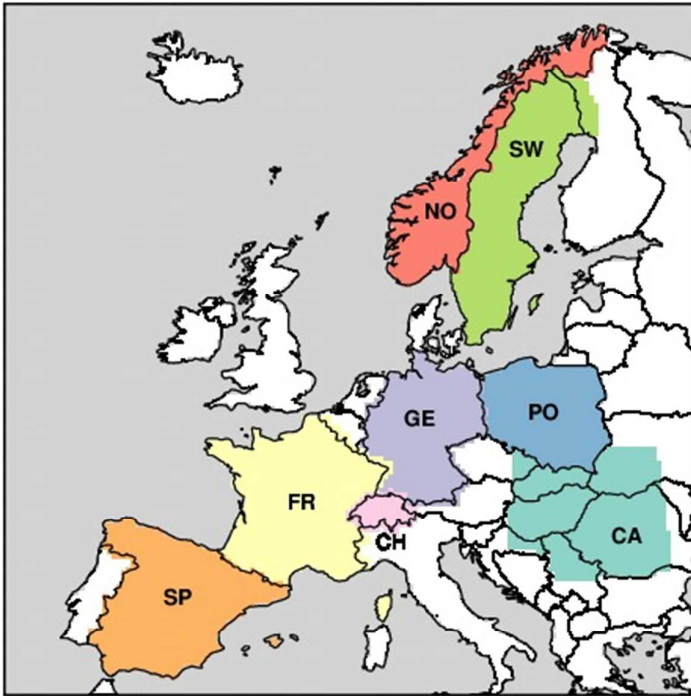
- 879 Rupp DE, Abatzoglou JT, Hegewisch KC, Mote PW. 2013. Evaluation of CMIP5 20<sup>th</sup> century climate  
880 simulations for the Pacific Northwest USA. *Journal of Geophysical Research: Atmospheres* 118:  
881 10.884-10.906, doi:10.1002/jgrd.50843.
- 882 Räisänen J. 2007. How reliable are climate models? *Tellus A* 59: 2-29. doi: 10.1111/j.1600-  
883 0870.2006.00211.x.
- 884 Santer BD, Taylor KE, Gleckler PJ, Bonfils C, Barnett TP, Pierce DW, Wigley TML, Mears C, Wentz FJ,  
885 Brüggemann W, Gillett NP, Klein SA, Solomon S, Stott PA, Wehner MF. 2009. Incorporating model  
886 quality information in climate change detection and attribution studies. *Proceedings of the National  
887 Academy of Sciences* 106 (35): 14778-14783. doi: 10.1073/pnas.0901736106.
- 888 Schneider U, Becker A, Finger P, Meyer-Christoffer A, Ziese M, Rudolf B. 2014. GPCP's new land  
889 surface precipitation climatology based on quality-controlled in situ data and its role in quantifying  
890 the global water cycle. *Theoretical and Applied Climatology* 115: 15-40. doi: 10.1007/s00704-013-  
891 0860-x.
- 892 Sevruck B. 1985. Systematischer Niederschlagsmessfehler in der Schweiz. In: *Der Niederschlag in der  
893 Schweiz, Beiträge zur Geologie der Schweiz - Hydrologie*, Vol. 31, pp. 65– 75, Bundesamt für Wasser  
894 und Geologie, Bern, Switzerland.
- 895 Sevruck B. 1986. Correction of precipitation measurements summary report. In: *Correction of  
896 precipitation measurements*. Züricher Geographische Schriften 23: 13–23.
- 897 Sunyer MA, Sorup HJD, Christensen OB, Madsen H, Rosbjerg D, Mikkelsen PS, Arnbjerg-Nielsen. 2013.  
898 On the importance of observational data properties when assessing regional climate model  
899 performance of extreme precipitation. *Hydrology and Earth System Sciences* 17: 4323-4337. doi:  
900 10.5194/hess-17-4323-2013.
- 901 Szentimrey T. 2004. Multiple Analysis of Series for Homogenization (MASH); Verification procedure  
902 for homogenized time series". *Proceedings of the Fourth Seminar for Homogenization and Quality  
903 Control in Climatological Databases*, Budapest, Hungary. WMO, WCDMP-No. 56, pp. 193-201.
- 904 Szentimrey T, Bihari Z. 2007. Mathematical background of the spatial interpolation methods and the  
905 software MISH (Meteorological Interpolation based on Surface Homogenized Data Basis).  
906 *Proceedings from the Conference on Spatial Interpolation in Climatology and Meteorology*,  
907 Budapest, Hungary, 2004, COST Action 719, COST Office, 17–27. Lubin D, Massom R. 2006. *Polar  
908 Remote Sensing*. Volume I: Atmosphere and Oceans. Praxis Publishing Ltd: Chichester, UK.
- 909 Szentimrey T. 2013. Multiple Analysis of Series for Homogenization (MASH v3.03). CARPATCLIM  
910 Deliverable D2.10. Annex 3 – Description of MASH and MISH algorithms. Hungarian Meteorological  
911 Service. Available from <http://www.carpatclim-eu.org/docs/mashmish/mashmish.pdf> (last access: 03  
912 April 2017).
- 913 Tanarhte M, Hadjinicolaou P, Lelieveld J. 2012. Intercomparison of temperature and precipitation  
914 data sets based on observations in the Mediterranean and the Middle East. *Journal of Geophysical  
915 Research* 117: D12102. doi: 10.1029/2011JD017293.
- 916 Uboldi F, Lussana C, Salvati M. 2008. Three-dimensional spatial interpolation of surface  
917 meteorological observations from high-resolution local networks. *Meteorological Applications* 15:  
918 331-345. doi:10.1002/met.76.
- 919 Vidal JP, Martin E, Franchistéguy L, Baillon M, Soubeyroux JM. 2010. A 50-year high-resolution  
920 atmospheric reanalysis over France with the Safran system. *International Journal of Climatology* 30  
921 (11): 1627-1644.



- 922 Wagner S, Kunstmann H, Bardossy A. 2007. Uncertainties in water balance estimations due to scarce  
923 meteorological information: a case study for the White Volta catchment in West Africa. In:  
924 Quantification and Reduction of Predictive Uncertainty for Sustainable Water Resources  
925 Management (Proceedings of Symposium HS2004 at IUGG2007, Perugia, July 2007). IAHS Publ. 313.
- 926 Weigel AP, Knutti R, Liniger MA, Appenzeller C. 2010. Risks of model weighting in multimodel climate  
927 projections. *Journal of Climate* 23(15): 4175-4191. doi: 10.1175/2010JCLI3594.1.
- 928 Wolff MA, Isaksen K, Petersen-Overleir A, Odemark K, Reitan T, Braekkan R. 2015. Derivation of a  
929 new continuous adjustment function for correcting wind-induced loss of solid precipitation: results of  
930 a Norwegian field study. *Hydrology and Earth System Sciences* 19: 951-967. doi: 10.5194/hess-19-  
931 951-2015.
- 932

Peer Review Only

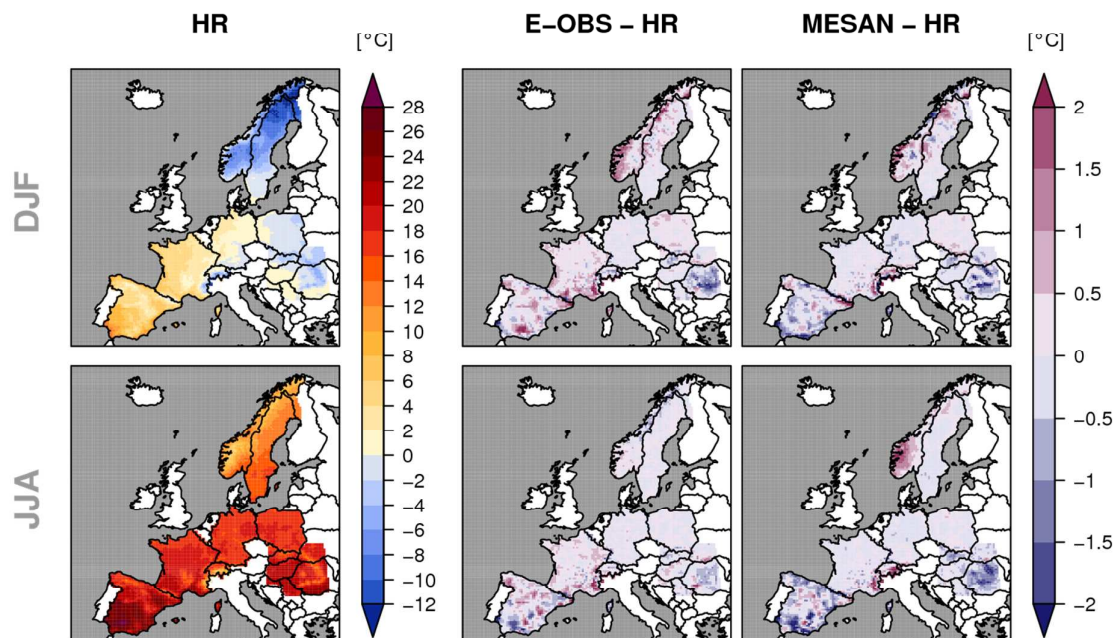
933 Figures



934

935 Fig. 1: The eight sub-regions considered for RCM evaluation. SP: Spain, FR: France, CH: Switzerland,  
936 GE: Germany, NO: Norway, SW: Sweden, PO: Poland, CA: Carpathians.

937



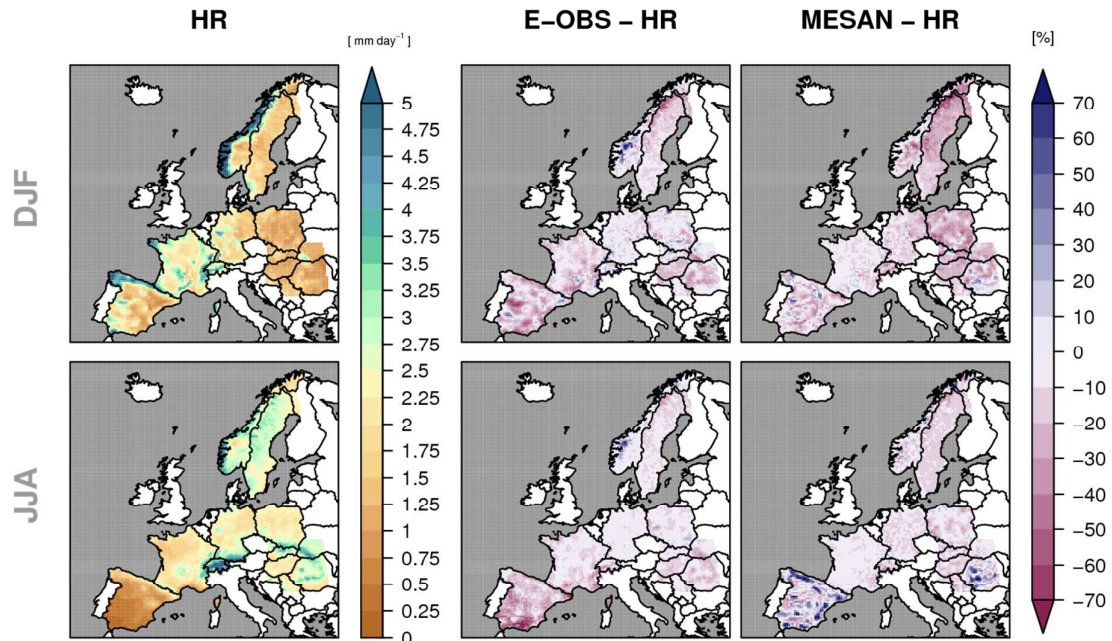
938

939

940

941

**Fig. 2:** Spatial pattern of seasonal mean temperature [ $^{\circ}\text{C}$ ] in HR in the period 1989-2006 (left column) and difference between E-OBS and HR (middle column) and MESAN and HR (right column). Upper row: Winter (DJF), lower row: Summer (JJA).



942

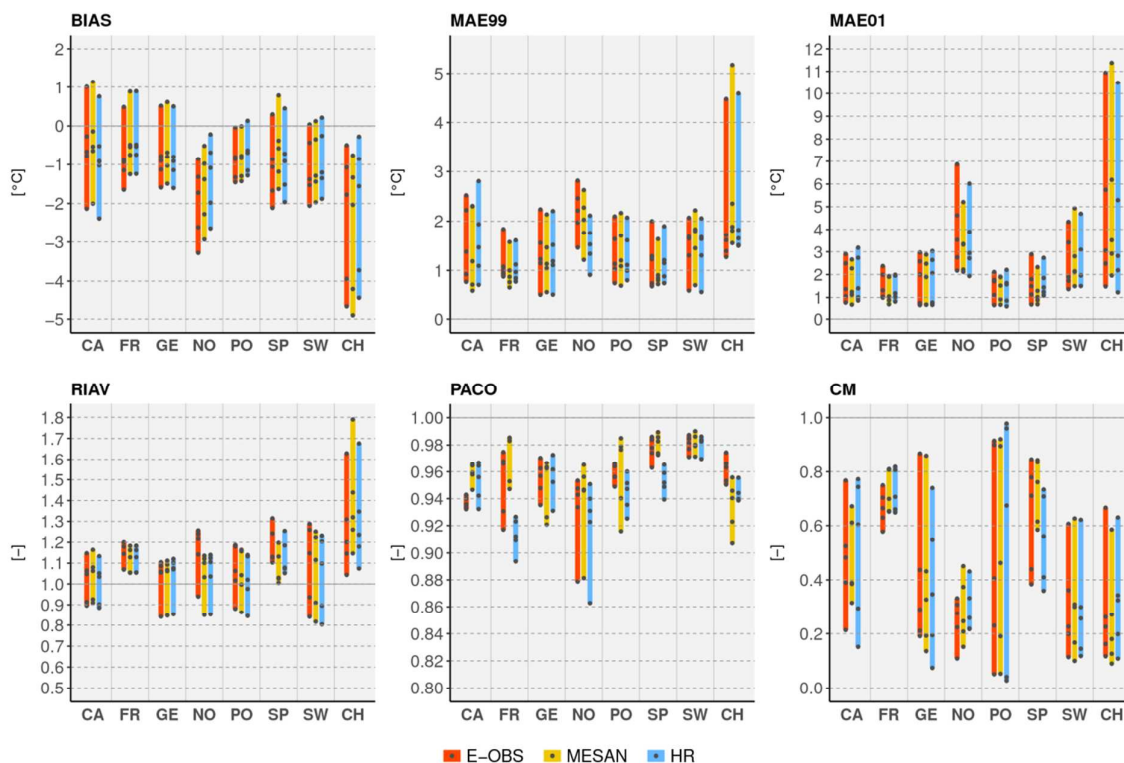
943

944

945

**Fig. 3:** As Figure 2 but for mean seasonal precipitation [ $\text{mm day}^{-1}$ ]. Differences between E-OBS and HR and between MESAN and HR are given in [%].

946



947

948

949

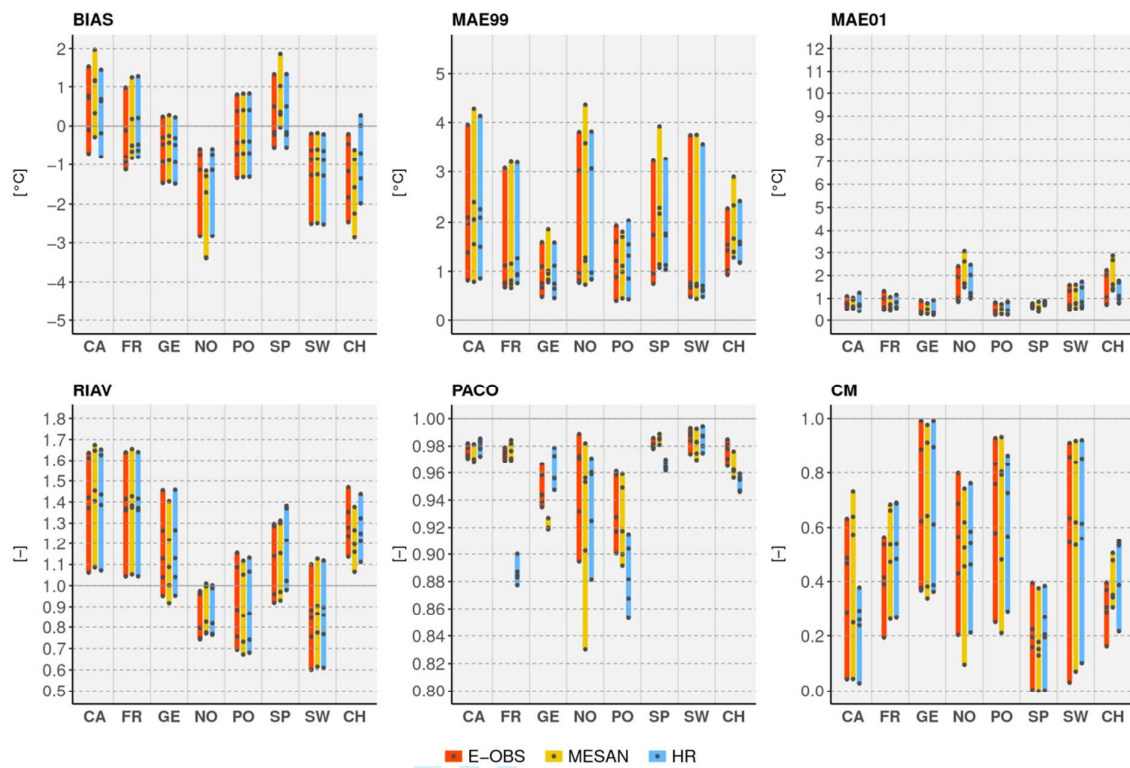
950

951

952

953

**Fig. 4:** Evaluation results for winter (DJF) temperature. The six panels correspond to the six performance metrics considered, the colours refer to the three observational references. Each set of three bars corresponds to one sub-region (x-axis). The five dots within each bar refer to the evaluation results for the five individual RCMs, whereas the bars themselves depict the model spread in terms of the minimum-maximum range.



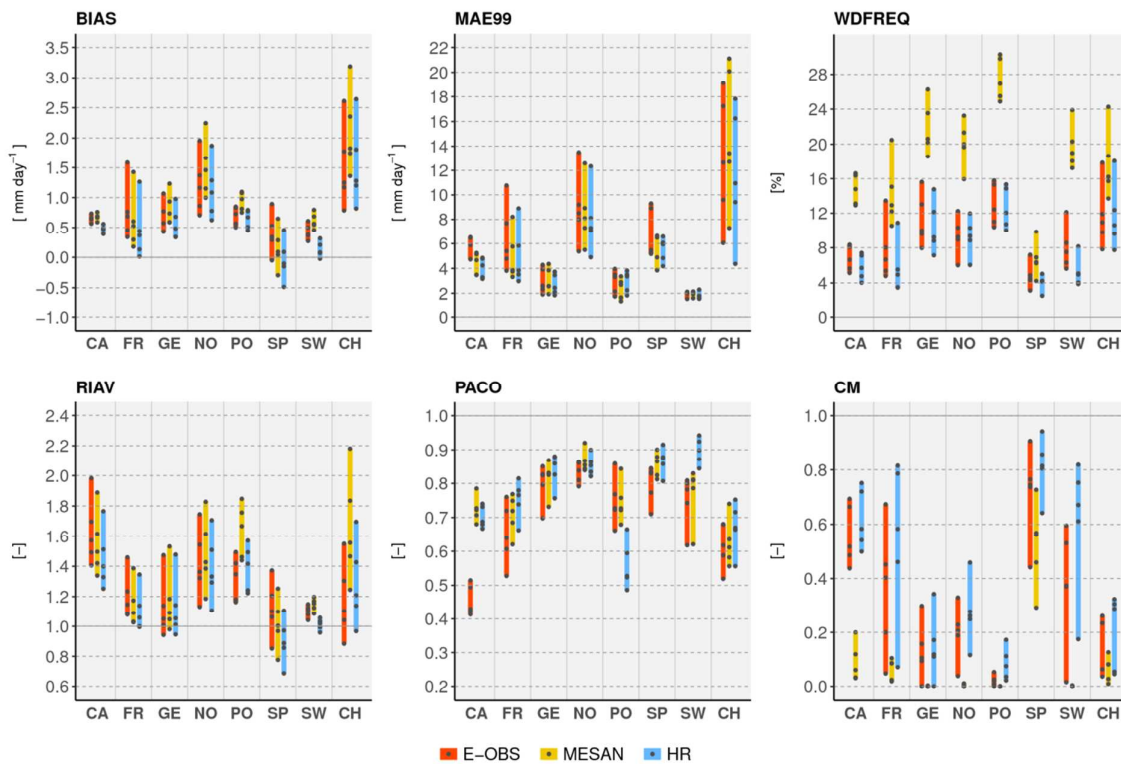
954

955

Fig. 5: As Figure 4 but for summer (JJA) temperature.

956

957



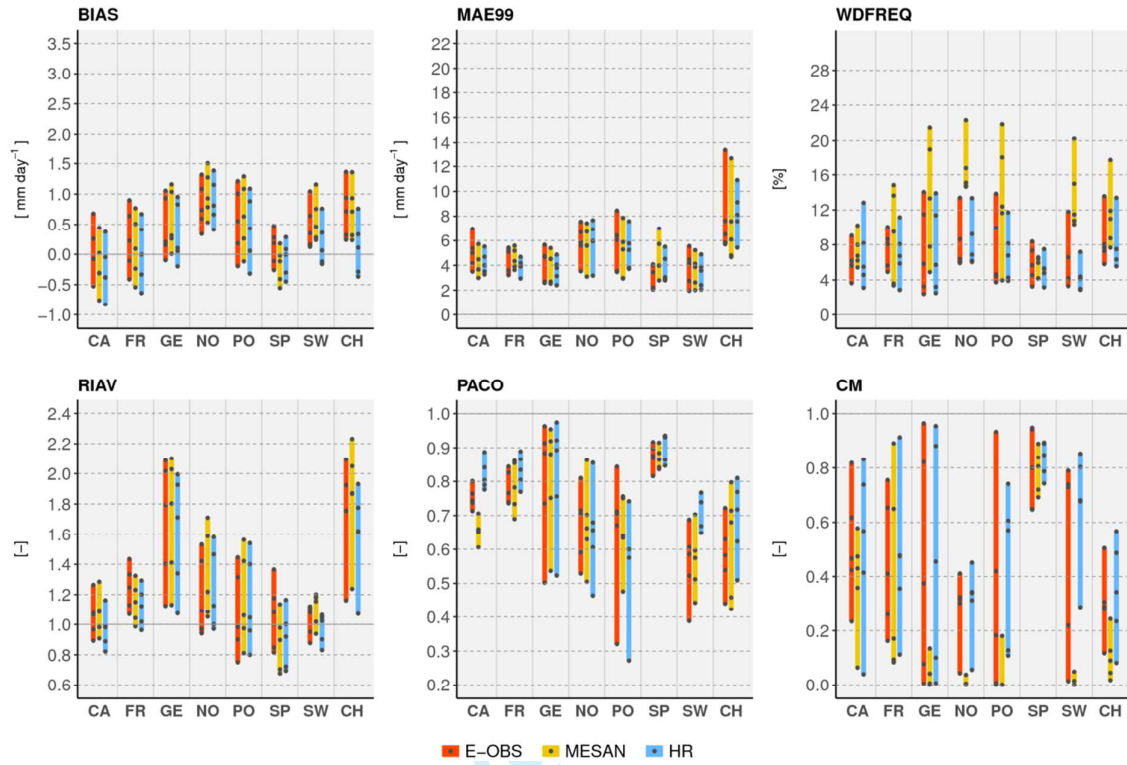
958

959 Fig. 6: As Figure 4 but for winter (DJF) precipitation.

960

961

962



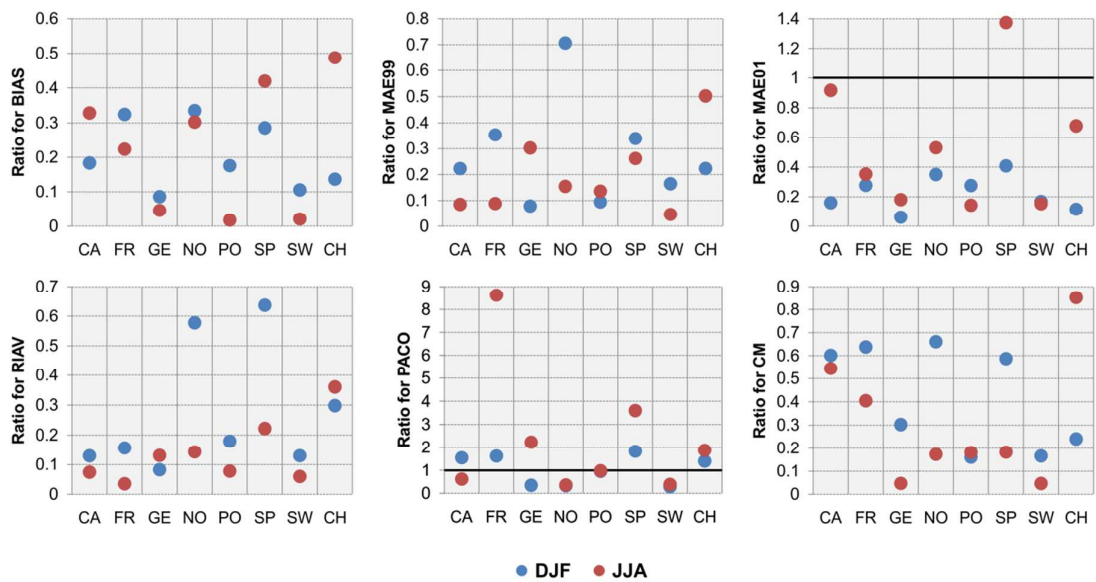
963

964 Fig. 7: As Figure 4 but for summer (JJA) precipitation.

965



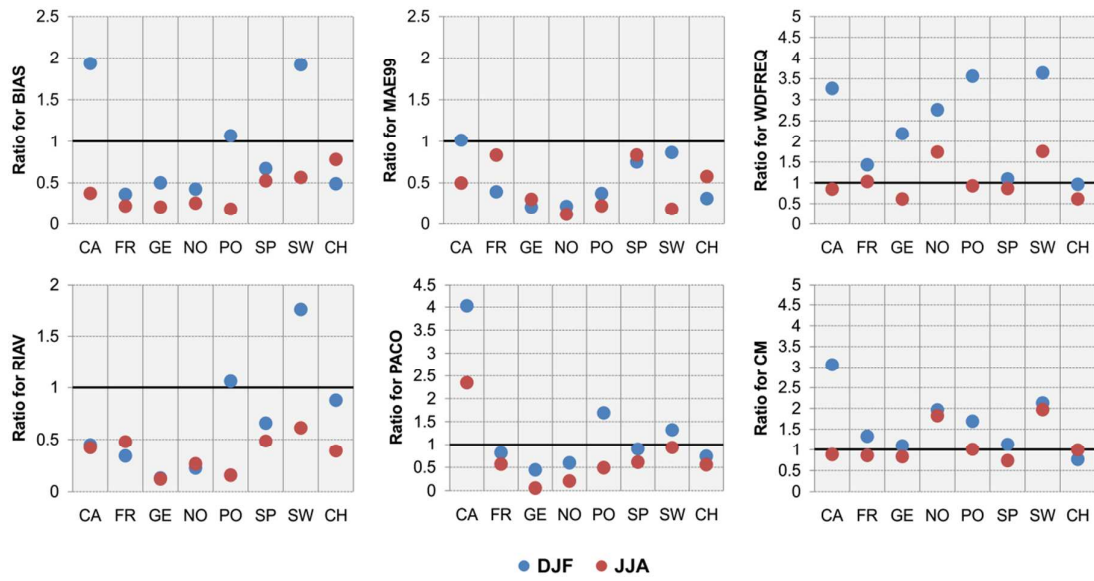
966



967 **Fig. 8:** Uncertainty intercomparison for temperature. The six panels refer to the six performance  
 968 metrics considered, the two colours to the seasons. An uncertainty ratio  $R$  larger (smaller)  
 969 (thick horizontal line) corresponds to a dominating observational (model) uncertainty for the  
 970 respective case.

971

972

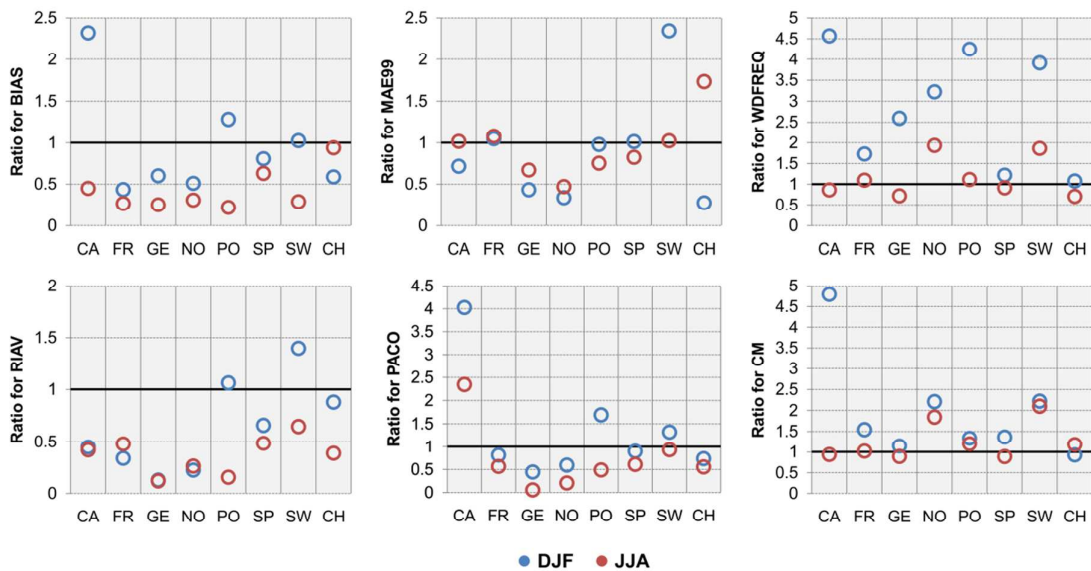


973

974 Fig. 9: As Figure 8 but for precipitation.

975

Peer Review Only

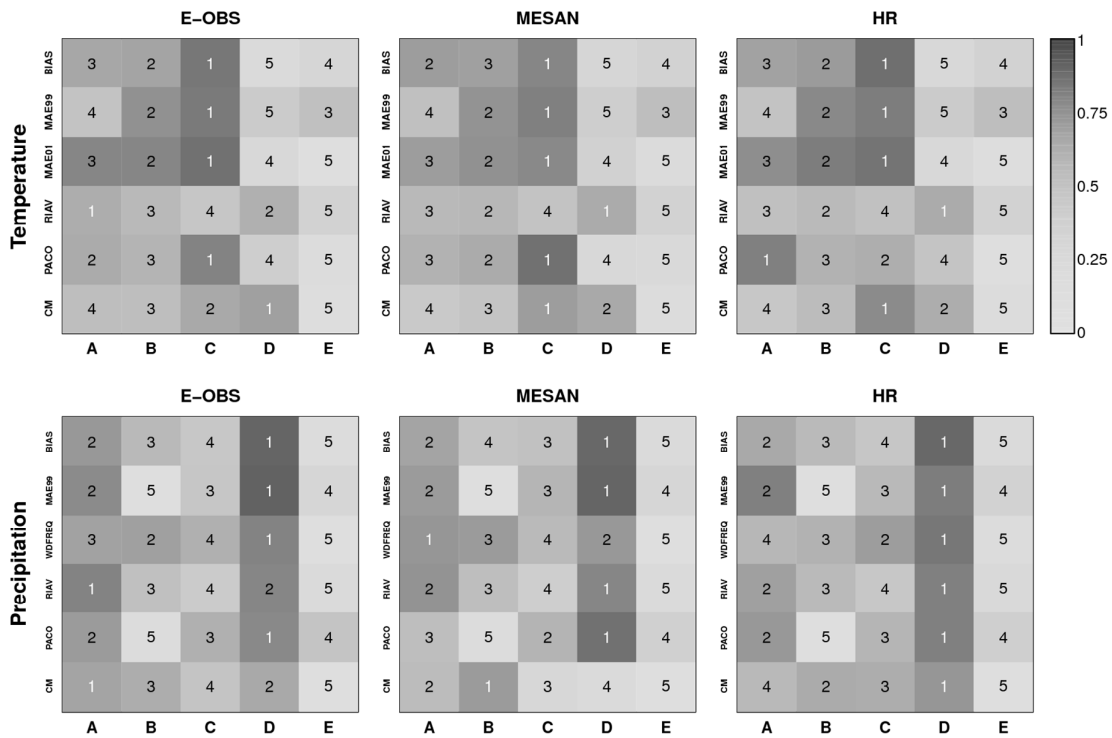


976

977 **Fig. 10:** As Figure 9 but for corrected precipitation: 20% were added to all daily precipitation amounts  
 978 in all three observational references except for HR over sub-region SW. Open circles instead of filled  
 979 ones are used for better separation from Fig. 9.

980

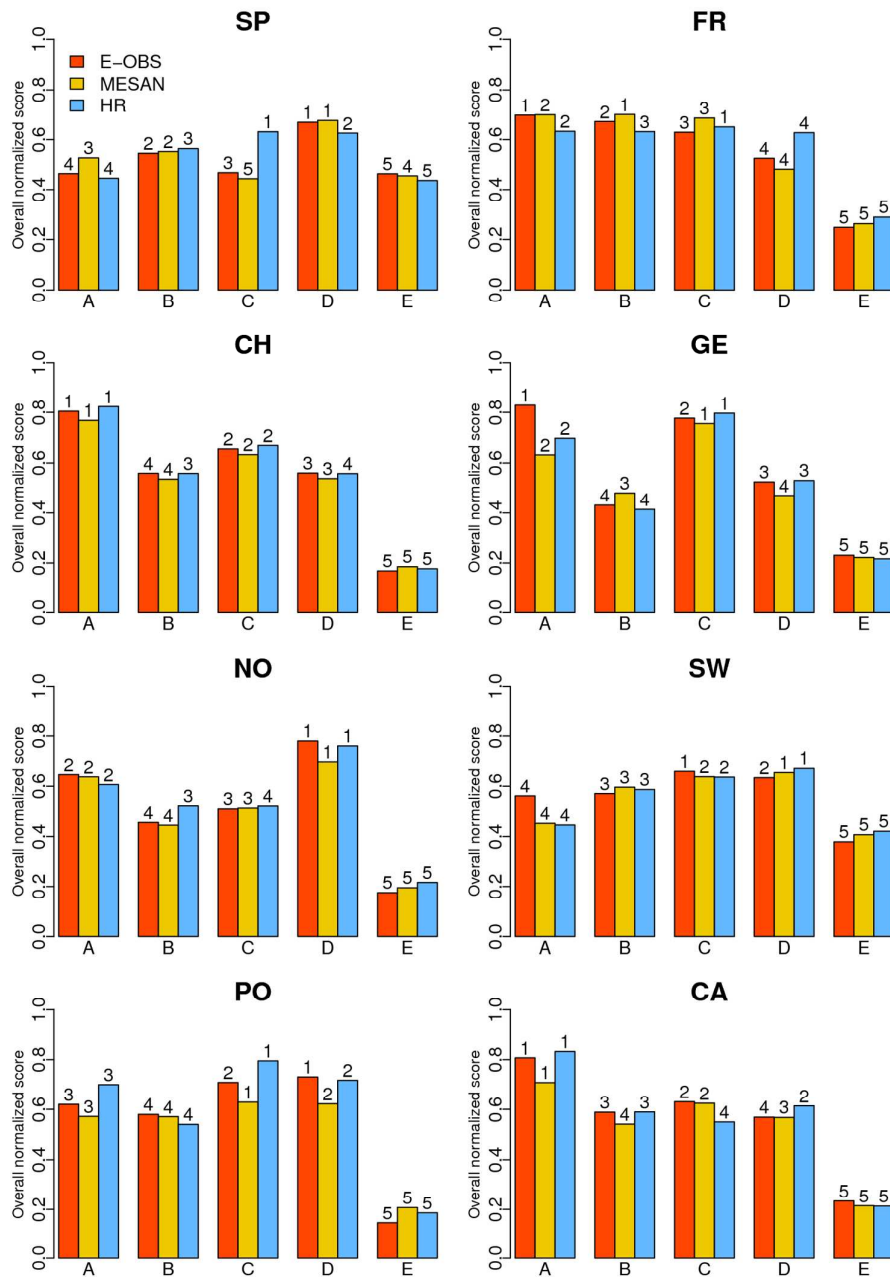
981



982

983 **Fig. 11:** Normalized performance scores (shading) for individual performance metrics, when averaged  
 984 over all seasons and regions. The upper row shows the results for temperature and the lower row for  
 985 precipitation. Numbering inside the shaded boxes indicates the actual RCM rank for each case. In  
 986 each panel, the individual rows indicate the performance metric, the individual columns the five  
 987 RCMs considered.

988



989

990 **Fig. 12:** Overall (combined temperature and precipitation) normalized performance scores for each  
 991 sub-region. The numbering above the bars indicates the actual RCM ranks separately for each  
 992 reference dataset.

993

994

995

996

997

998

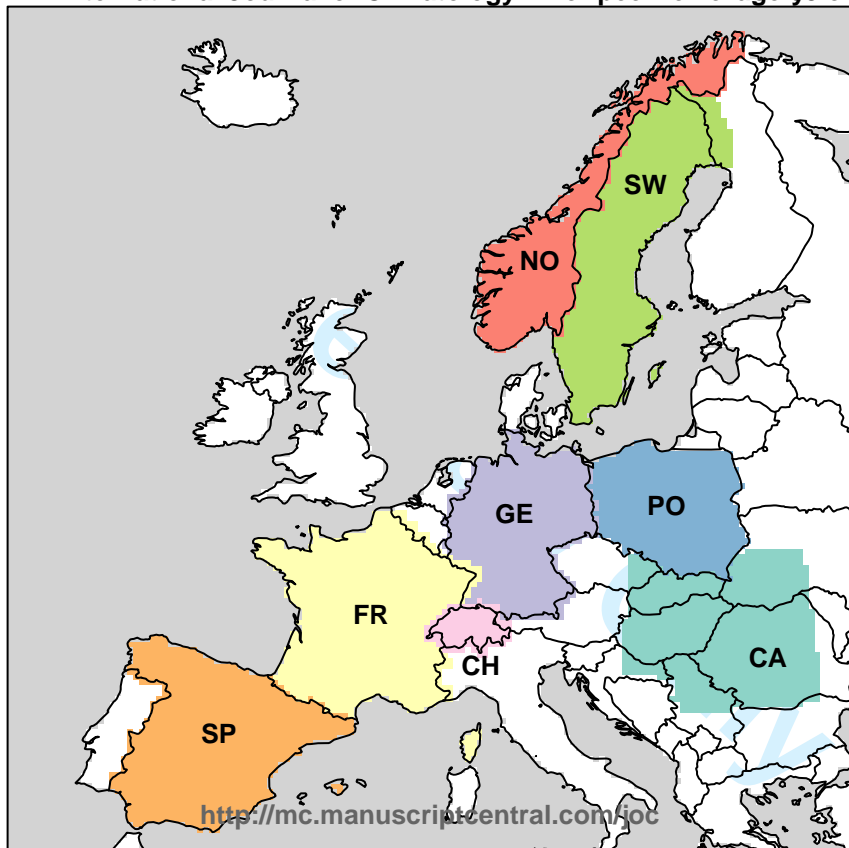
999

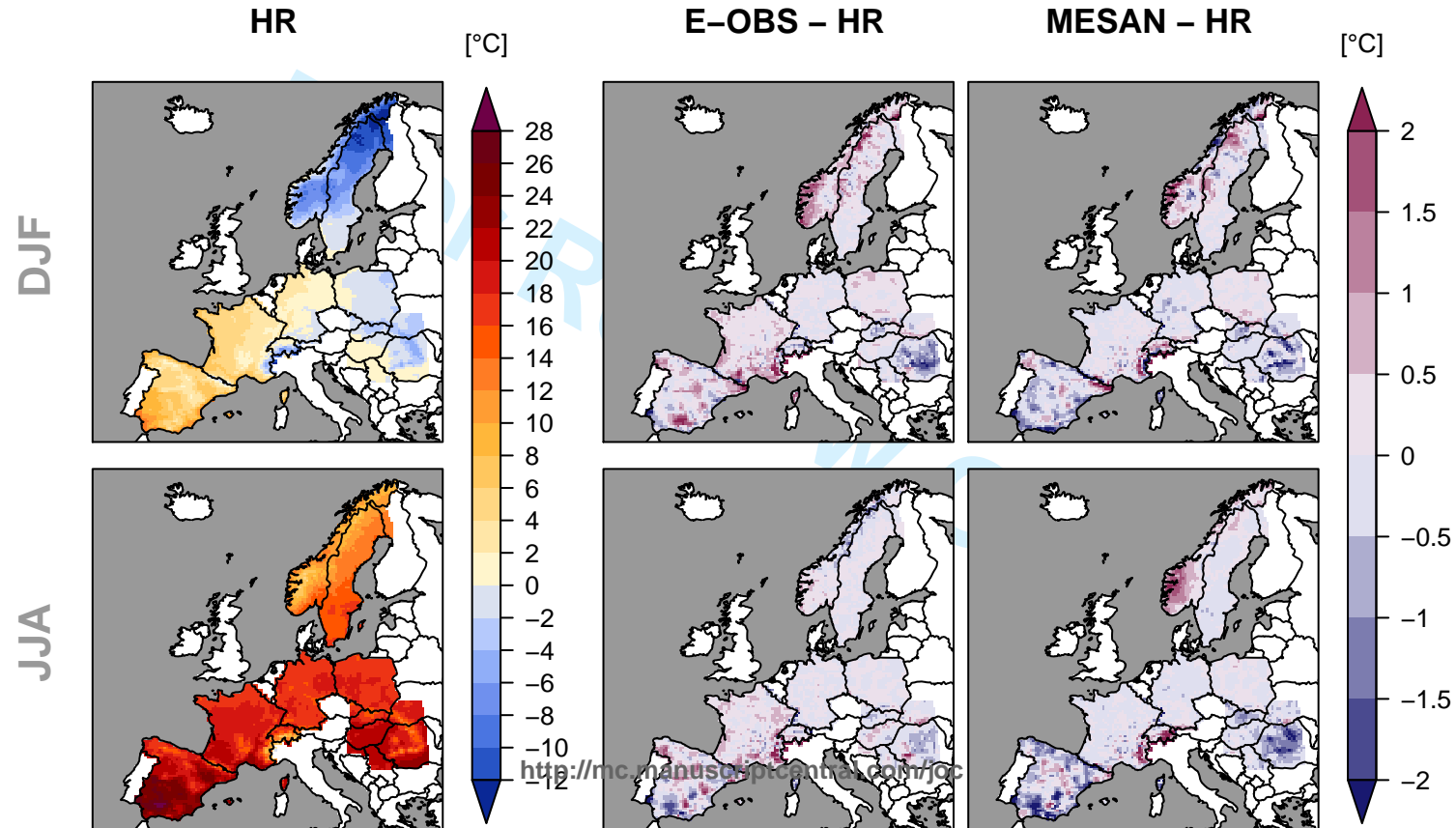
1000 **Tables**

1001 Table 1: Overview on the employed observational reference and RCM datasets. In this work the  
 1002 individual datasets are simply referred to by their abbreviation (last column).

<b>Type of dataset</b>	<b>Details</b>		
Observational reference	Name	Description	Abbreviation
	E-OBS v15	Section 2.1.1	<b>E-OBS</b>
	National high-resolution grids	Section 2.1.2	<b>HR</b>
RCM	EURO4M MESAN	Section 2.1.3	<b>MESAN</b>
	Model name and version	Institute/Group	Abbreviation
	CCLM 4.8.17	CLMcom	<b>A</b>
	HIRHAM 5	DMI	<b>B</b>
	WRF 3.3.1F	IPSL-INERIS	<b>C</b>
	RACMO 2.2E	KNMI	<b>D</b>
RCA 4	SMHI	<b>E</b>	

1003







HR

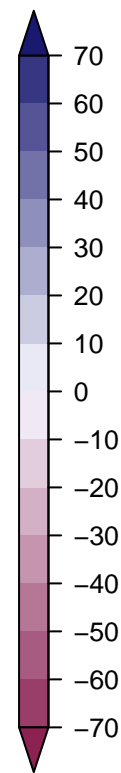
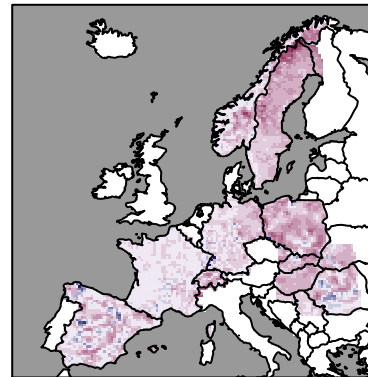
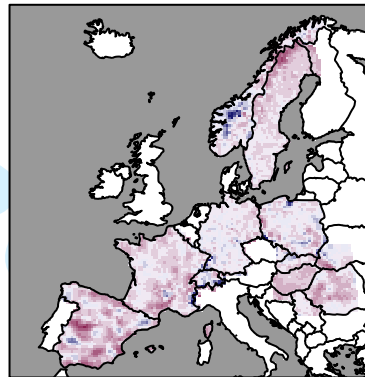
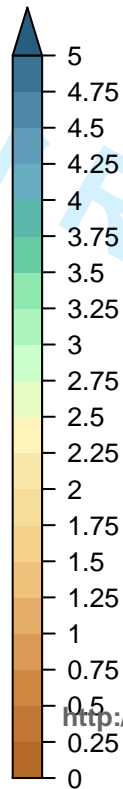
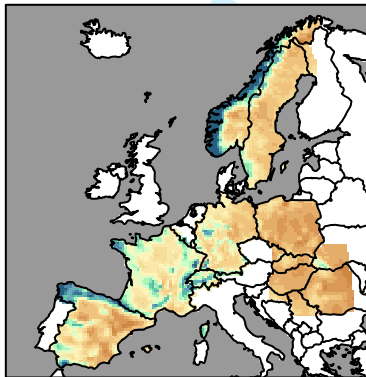
E-OBS – HR

MESAN – HR

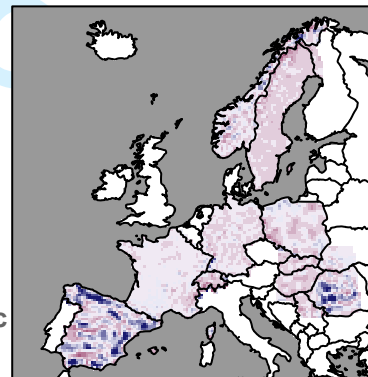
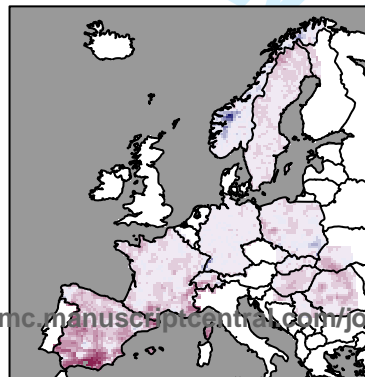
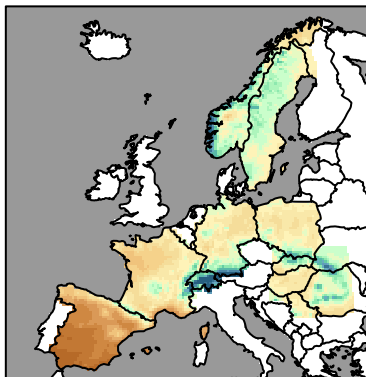
[ mm day<sup>-1</sup> ]

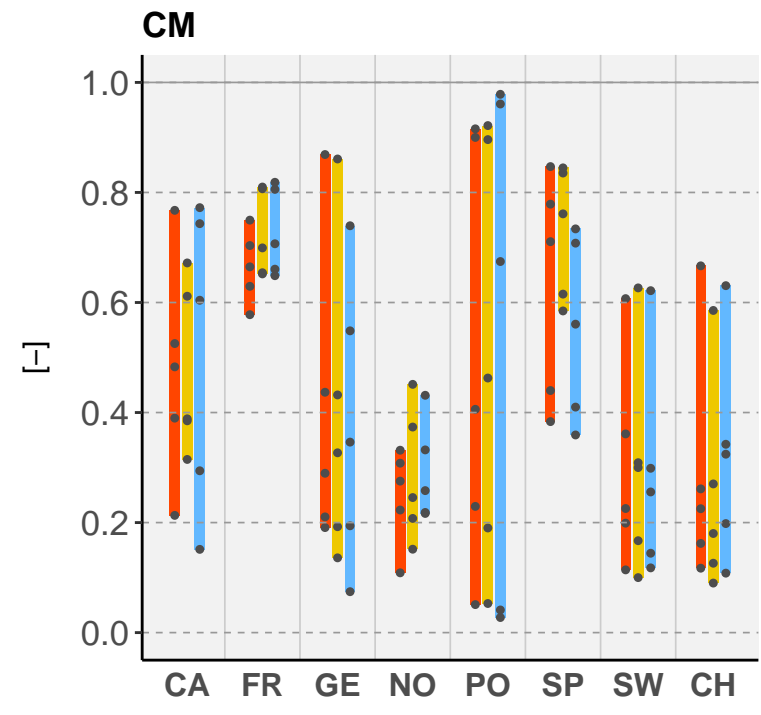
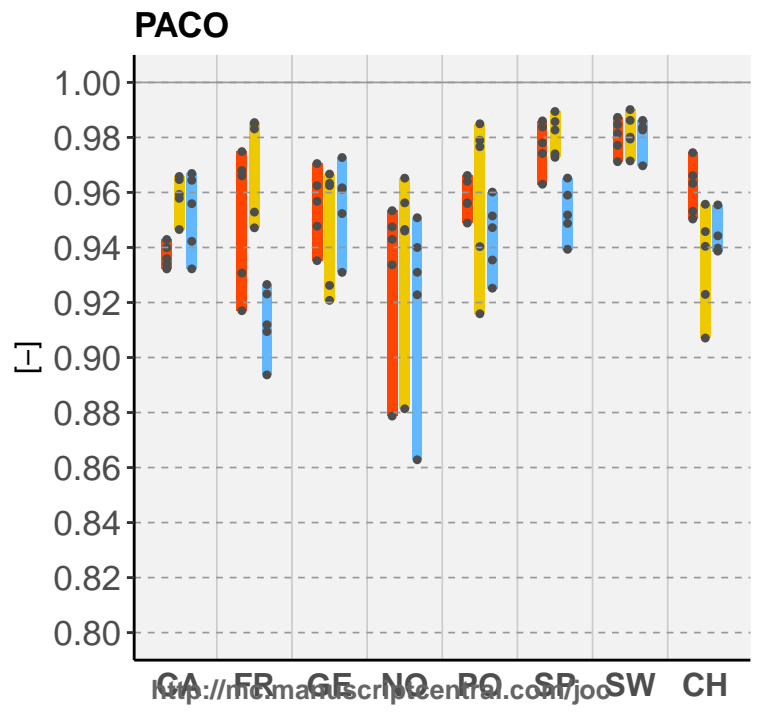
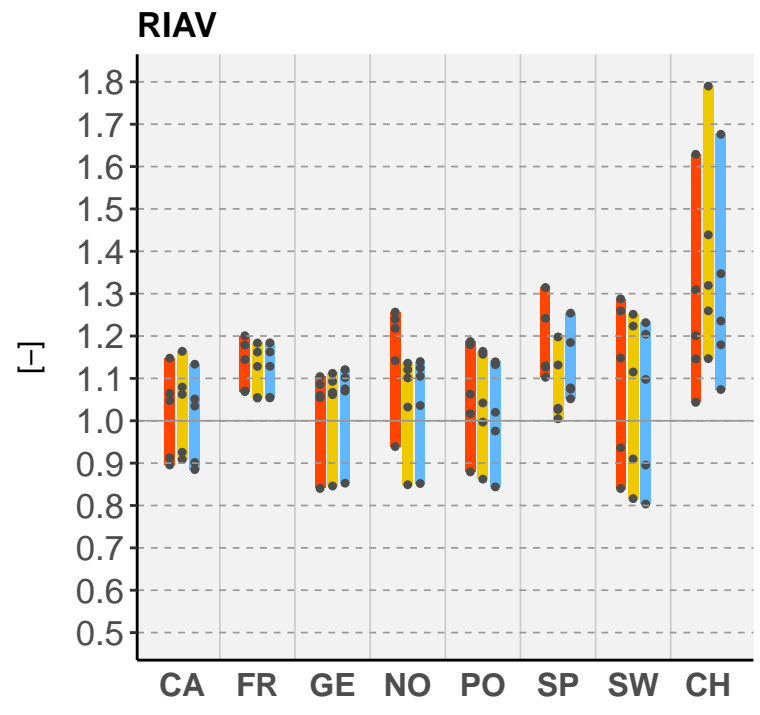
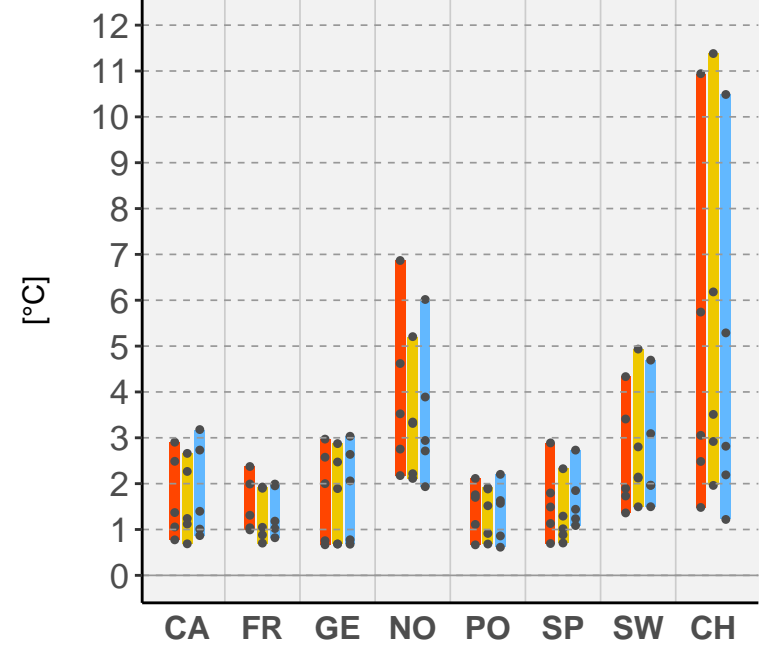
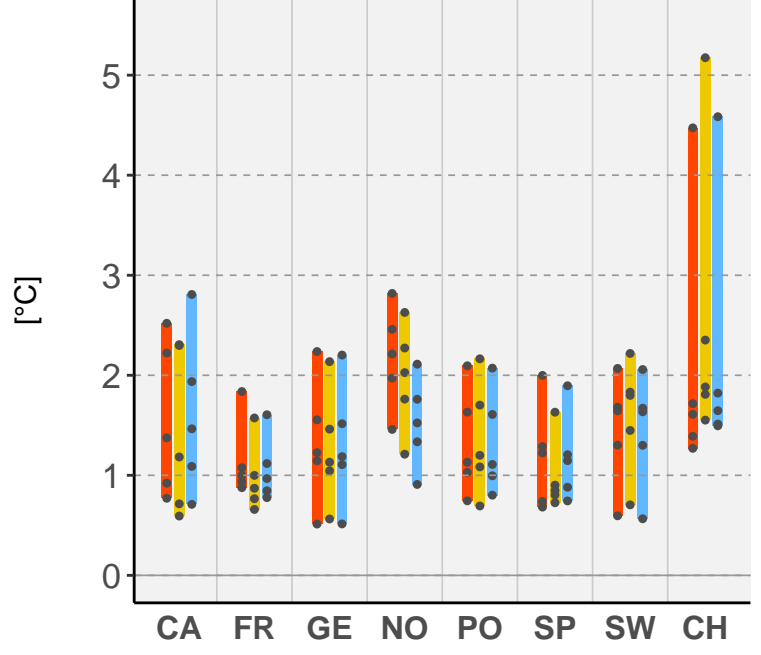
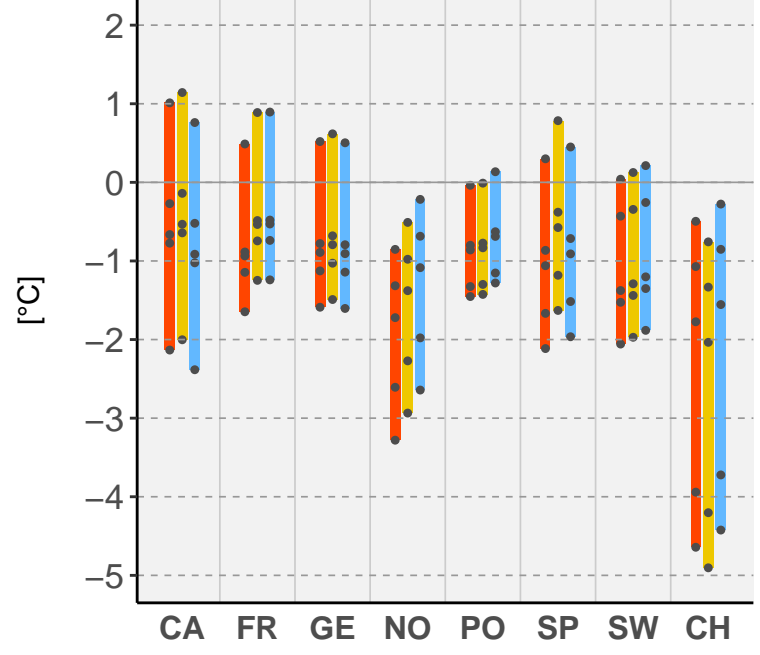
[%]

DJF

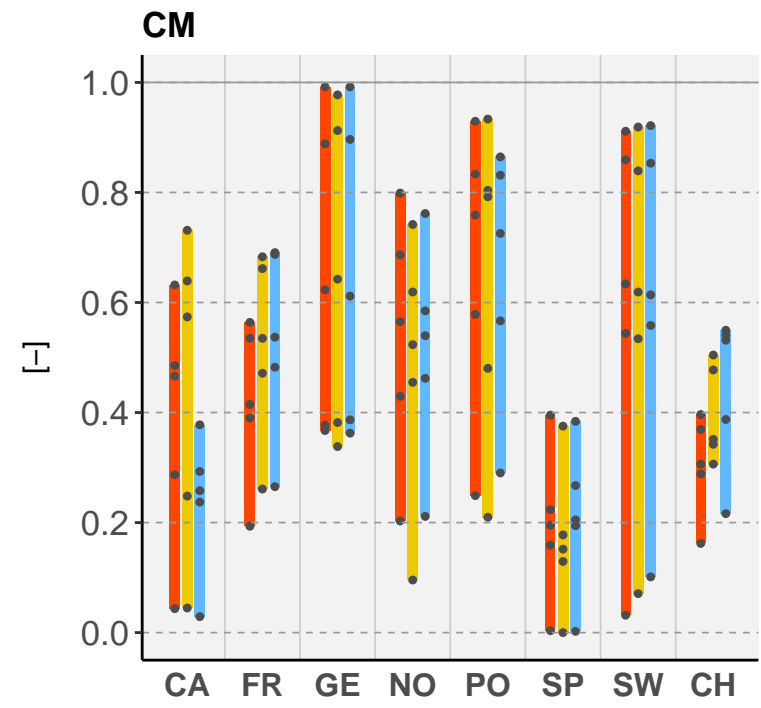
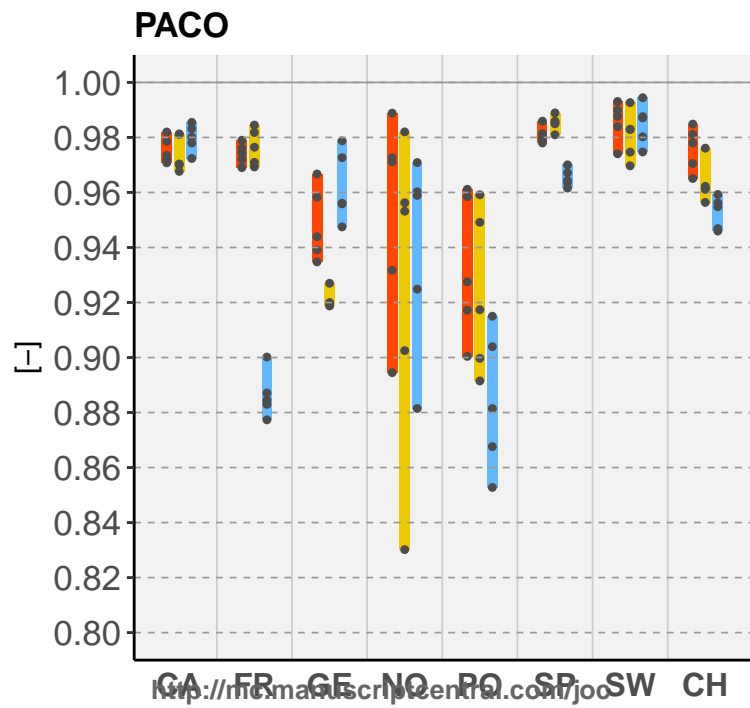
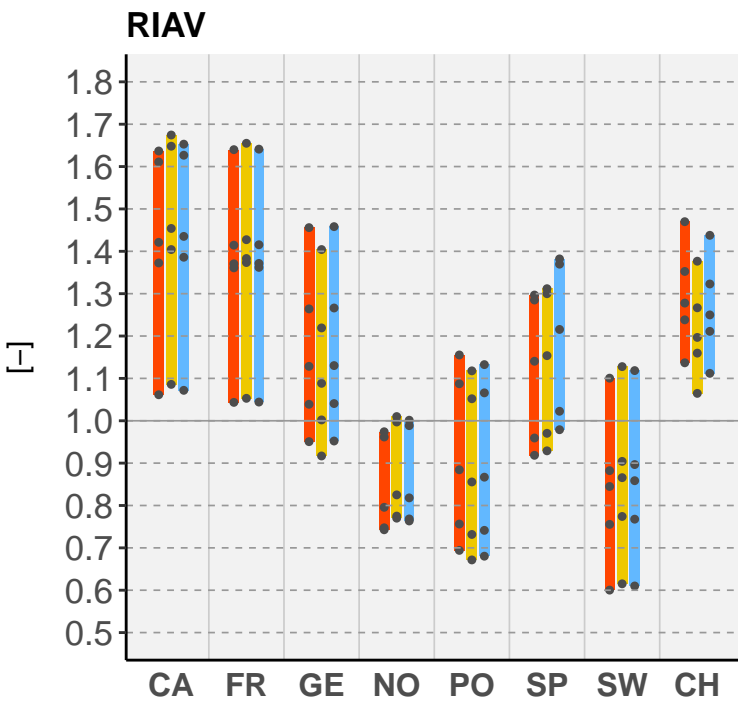
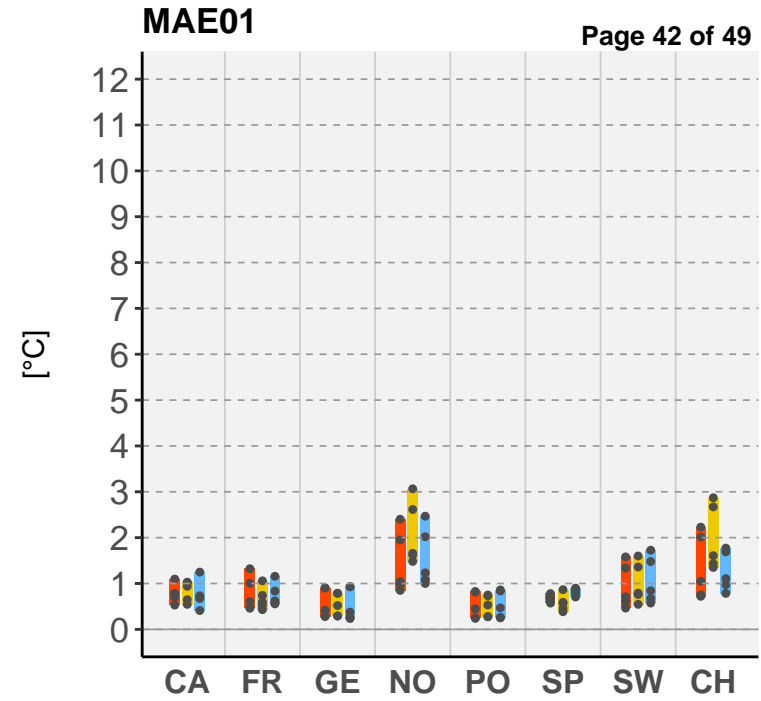
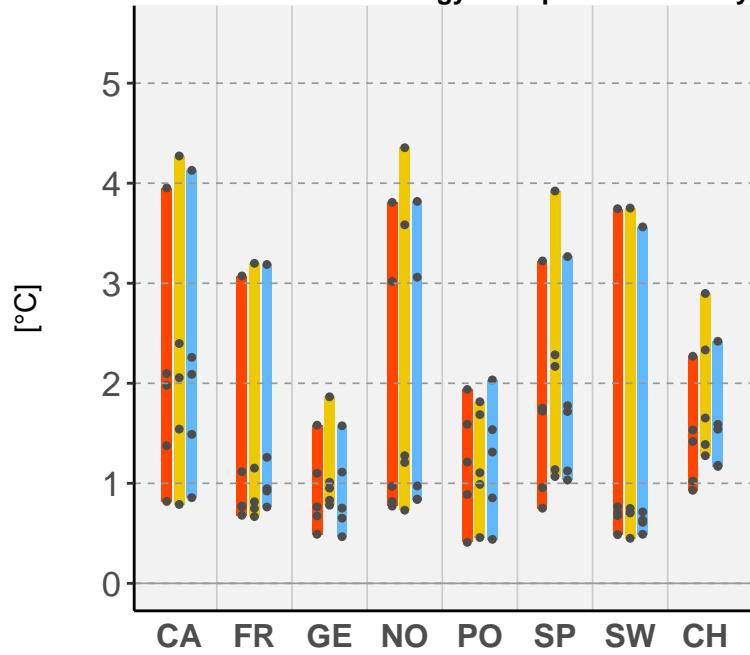
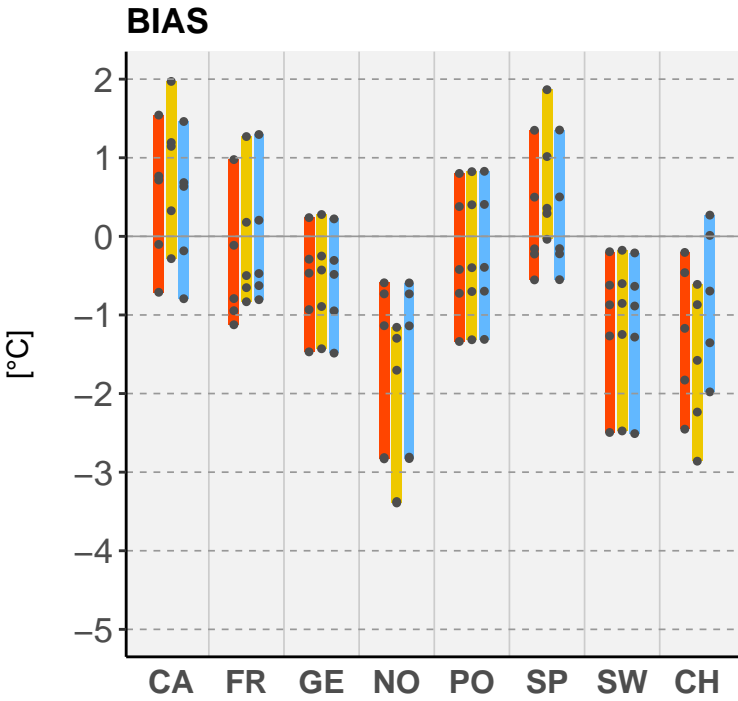


JJA

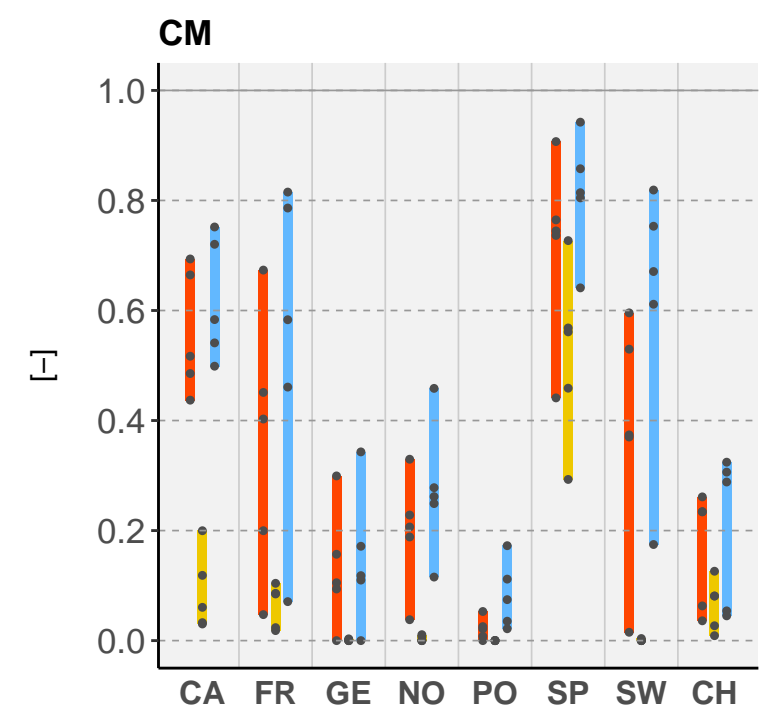
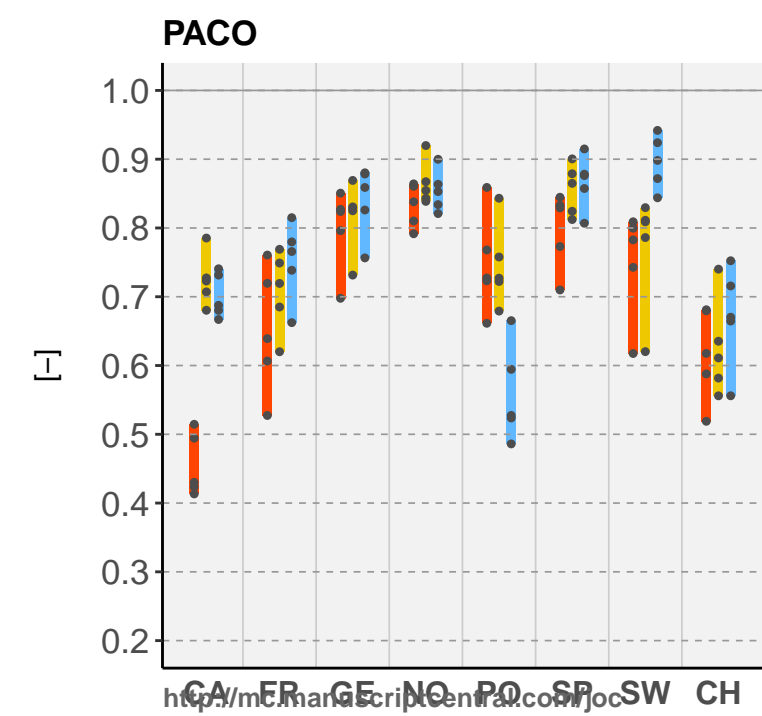
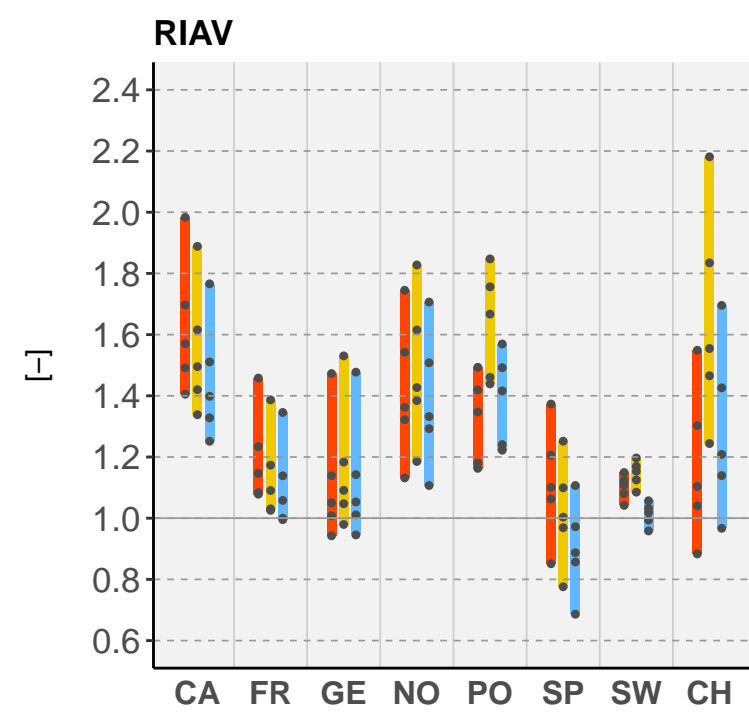
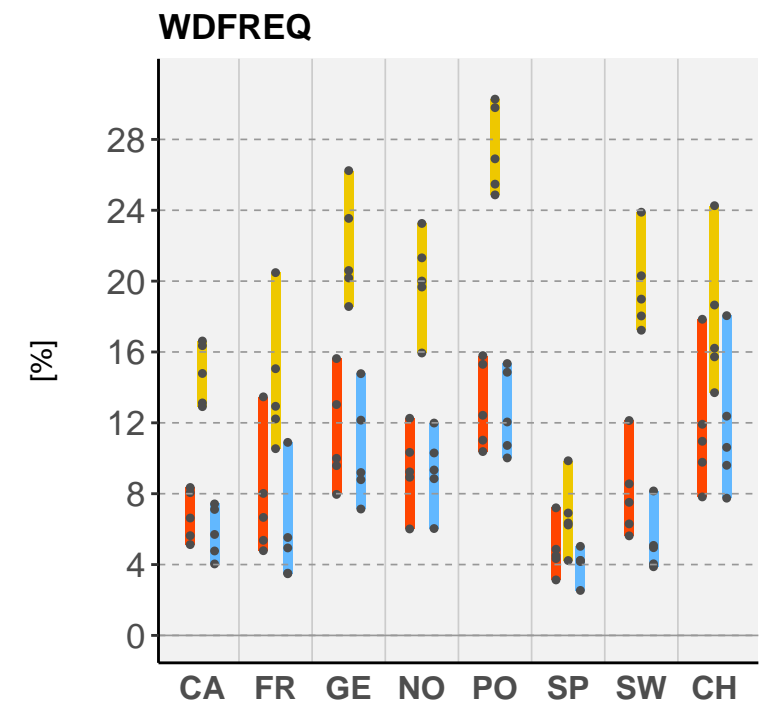
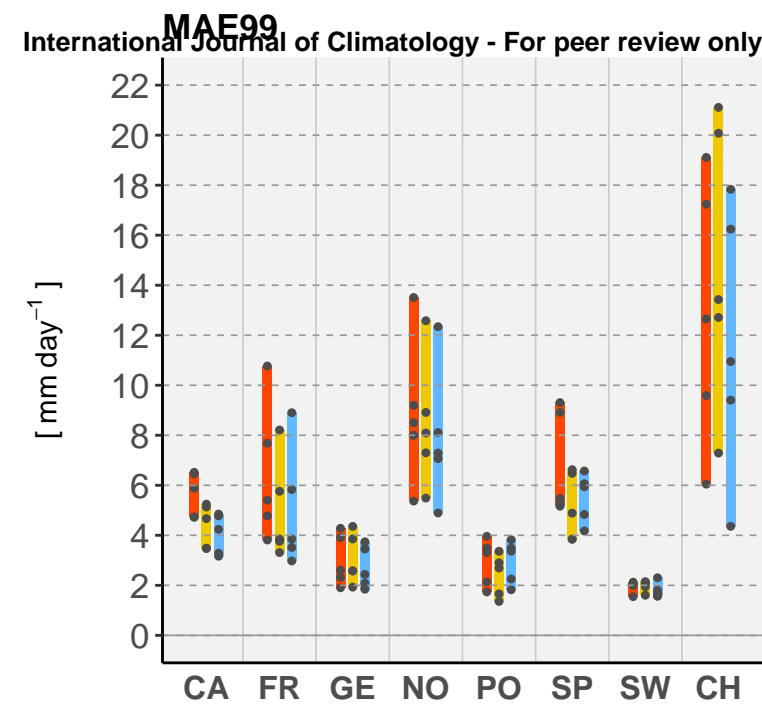
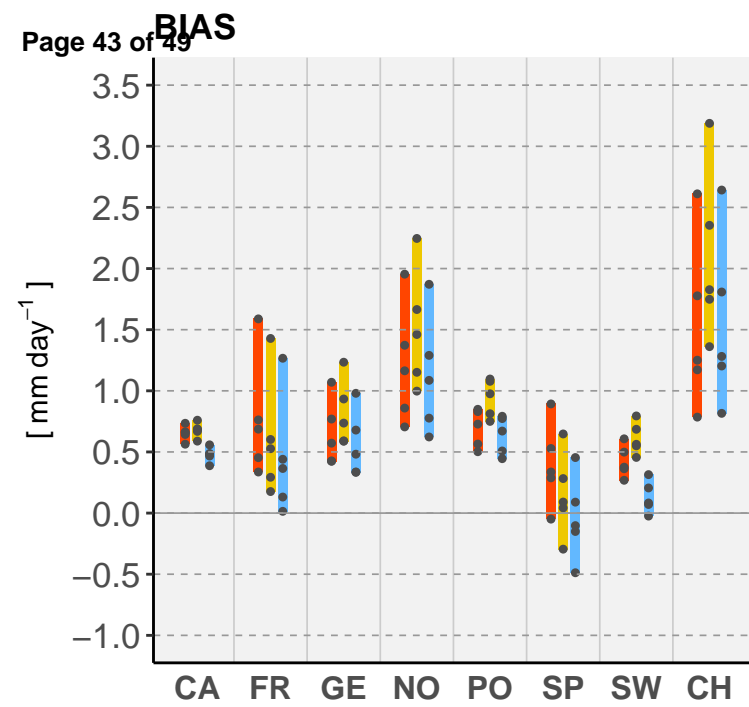




E-OBS MESAN HR

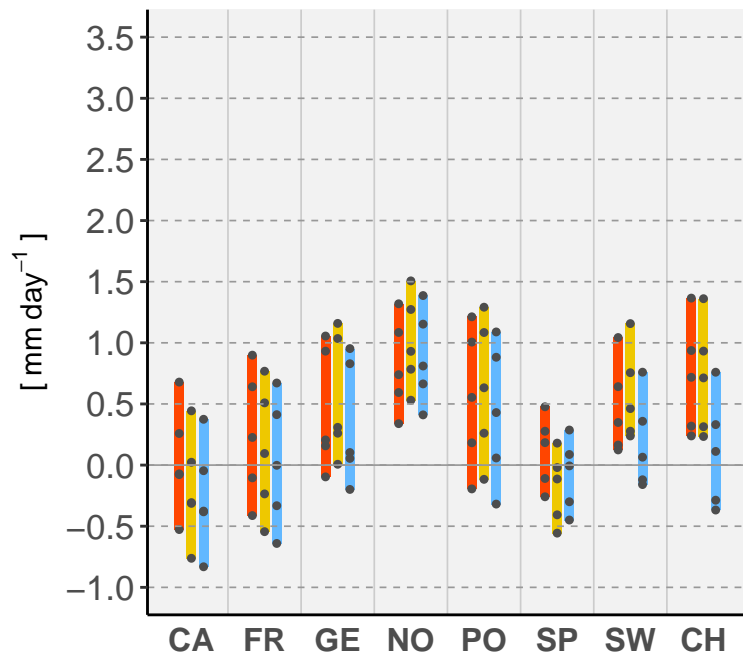


■ E-OBS   
 ■ MESAN   
 ■ HR

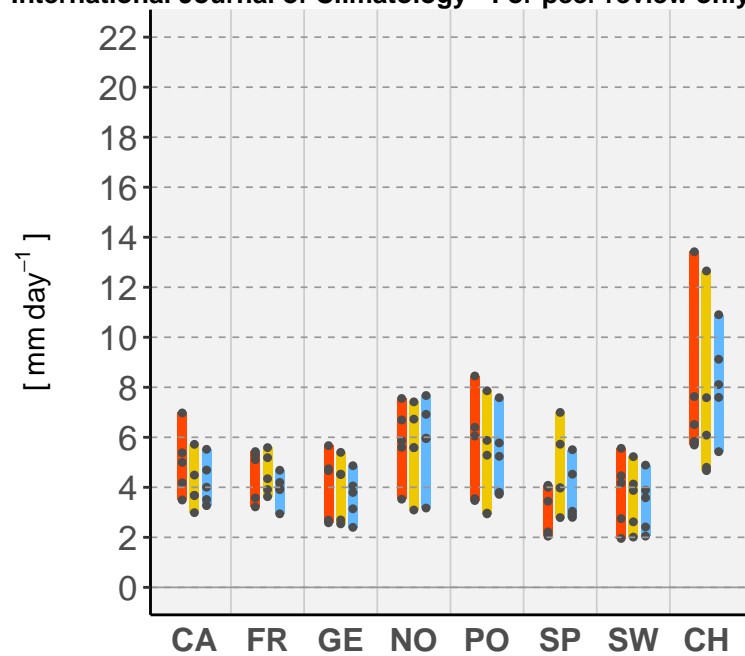


■ E-OBS   
 ■ MESAN   
 ■ HR

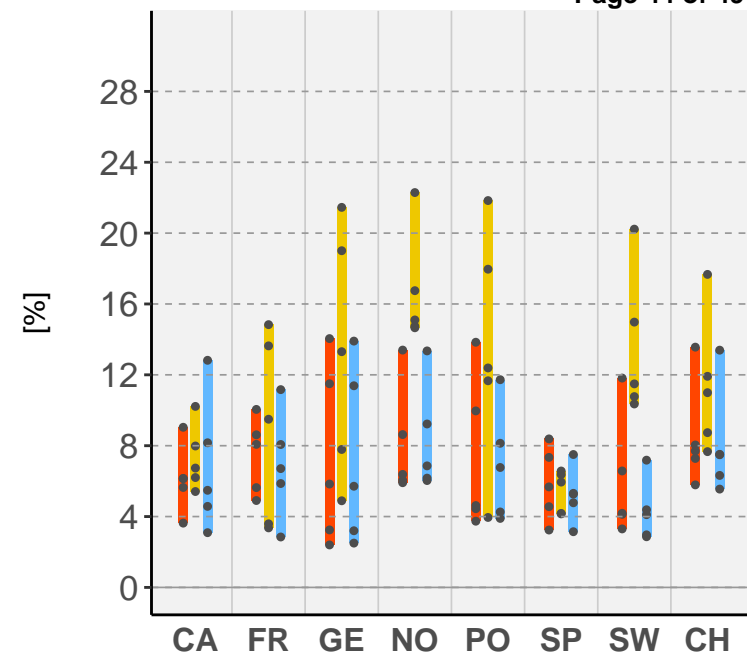
**BIAS**



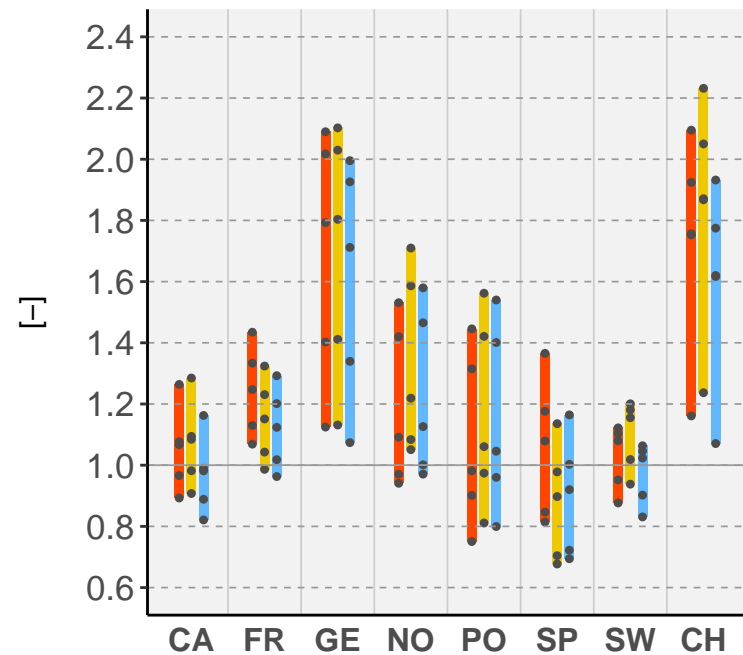
**MAE99**



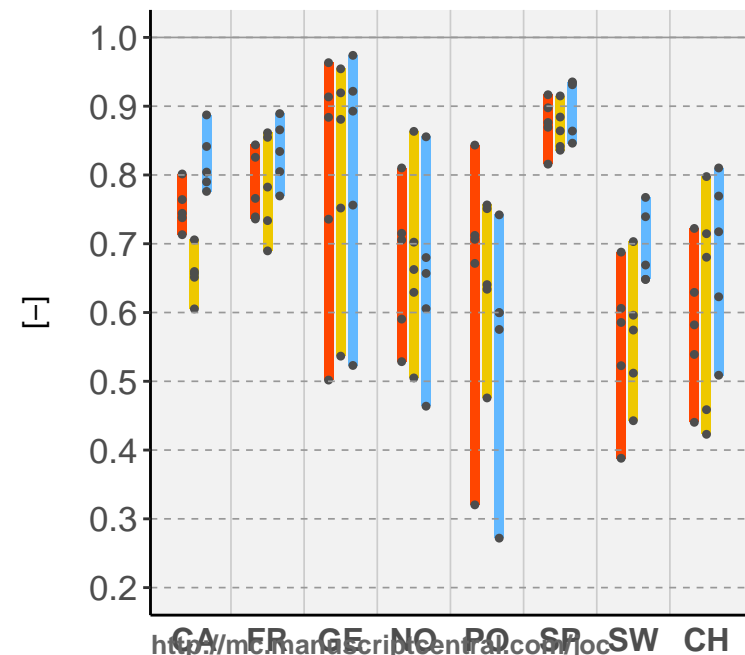
**WDFREQ**



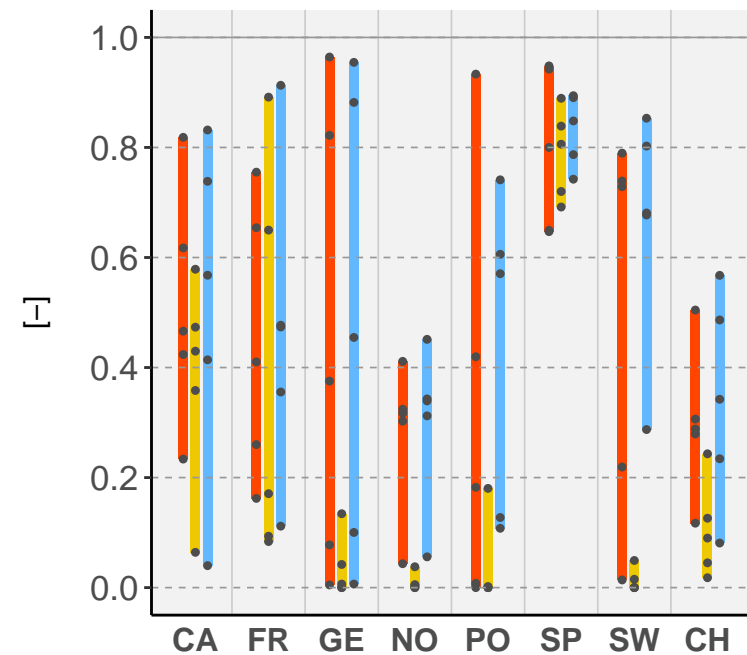
**RIAV**



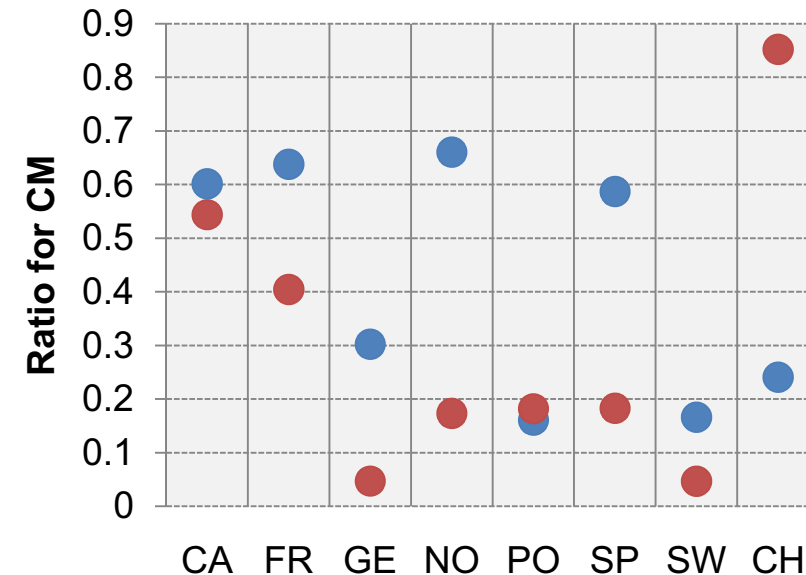
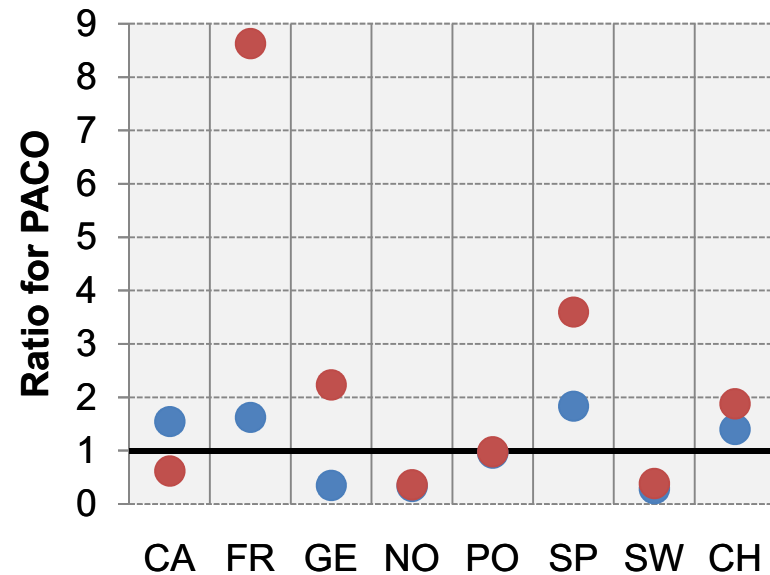
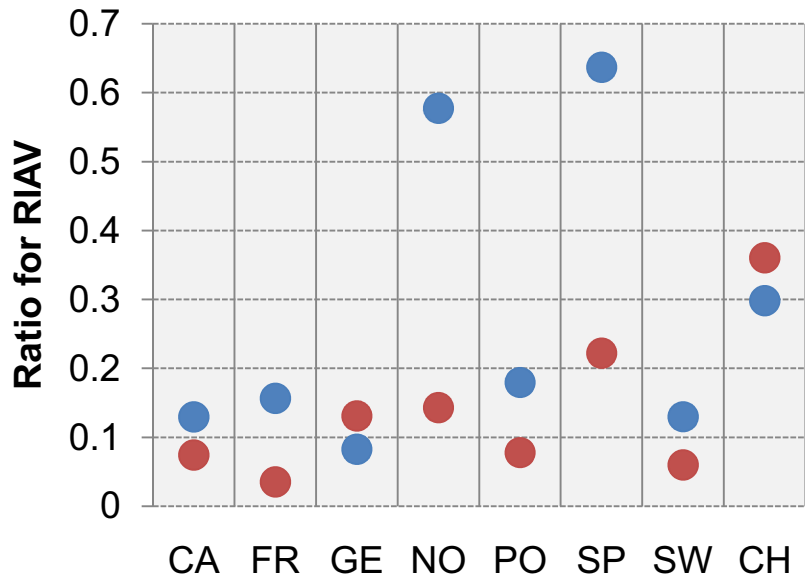
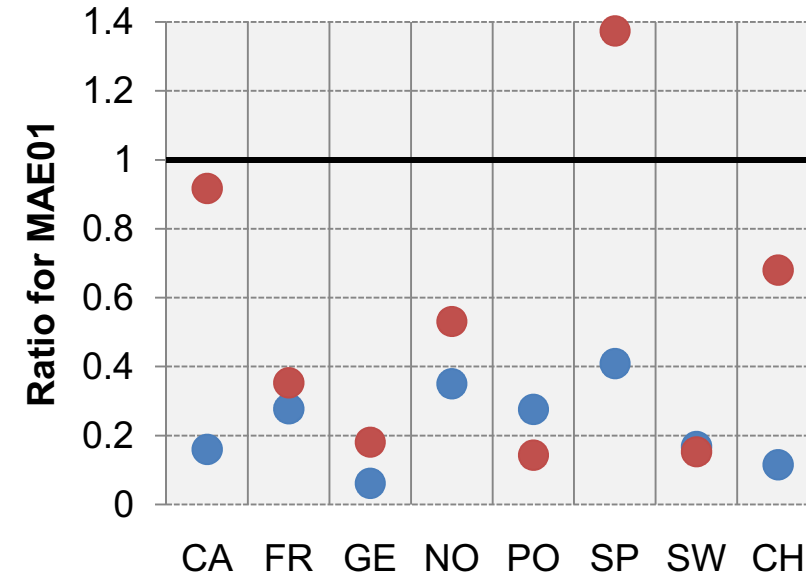
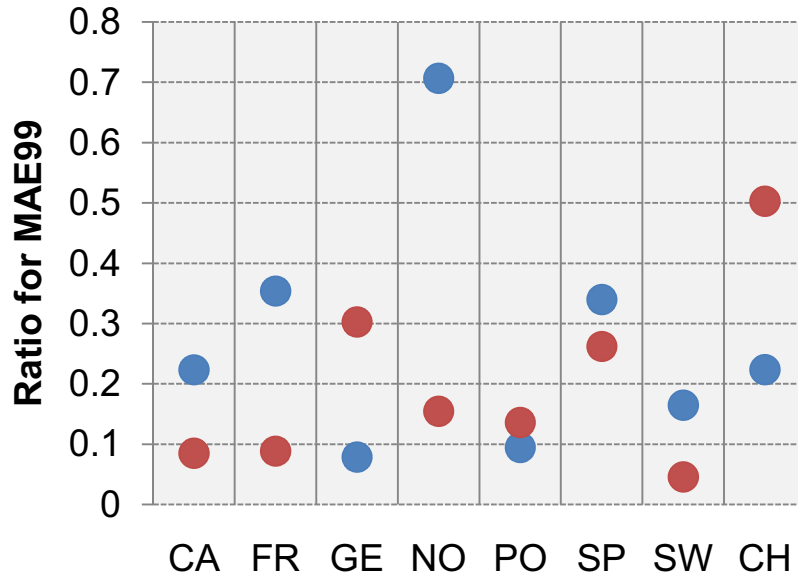
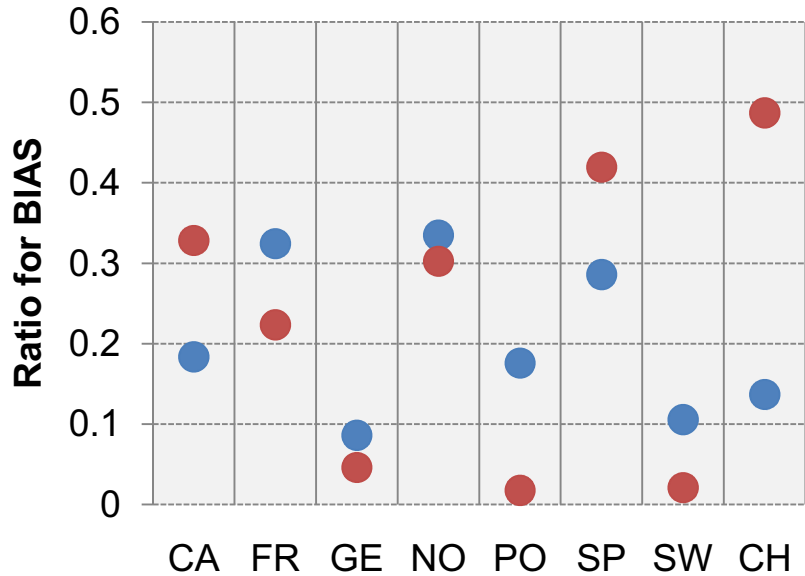
**PACO**

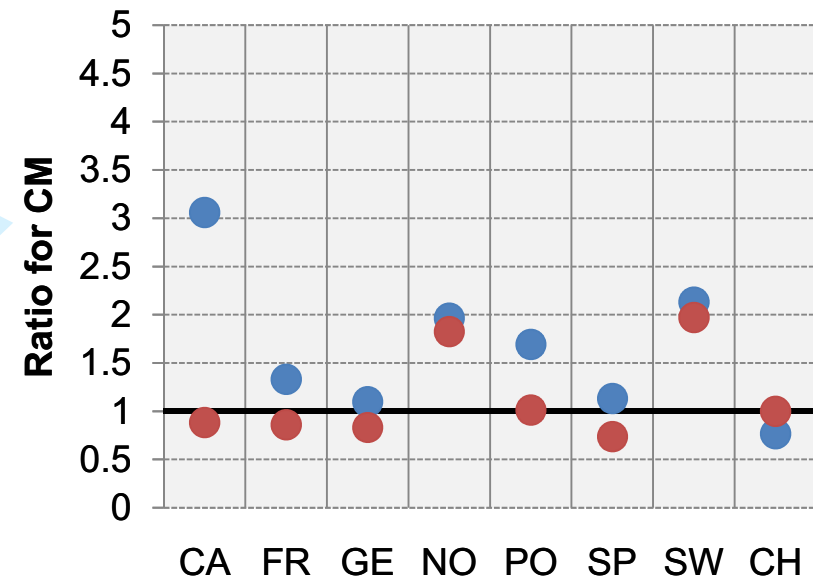
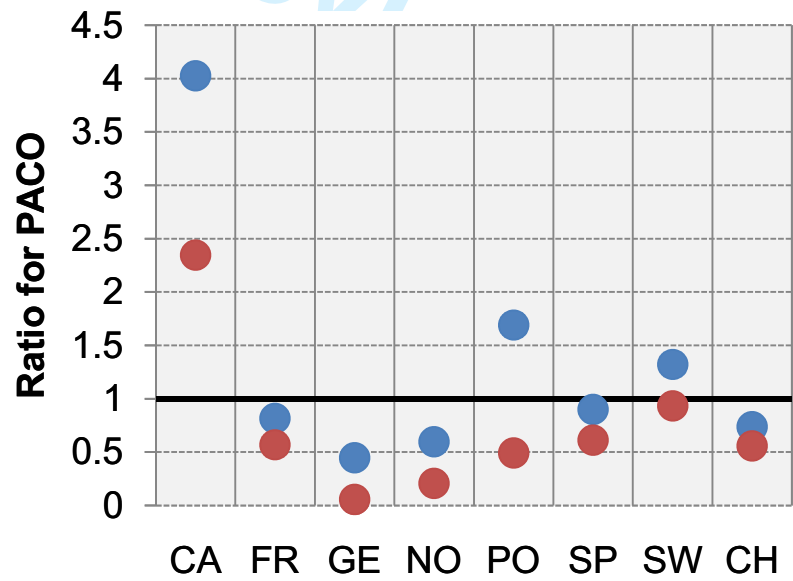
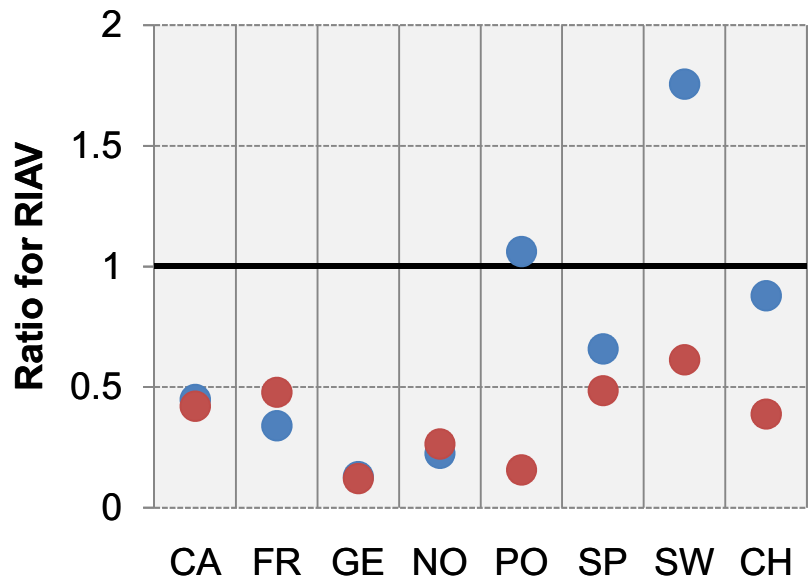
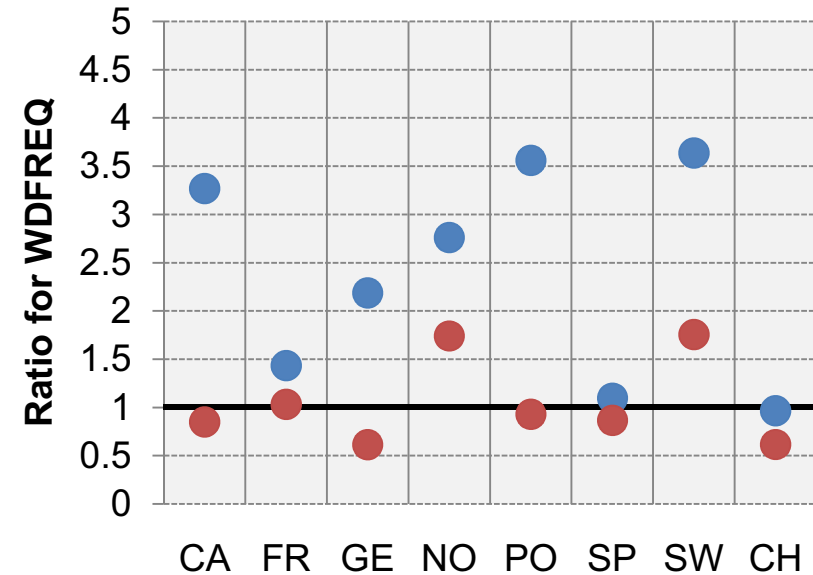
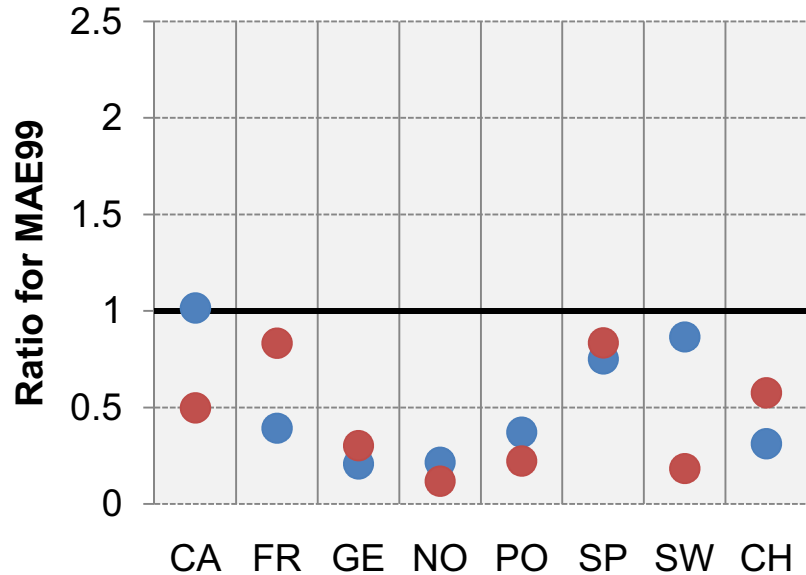
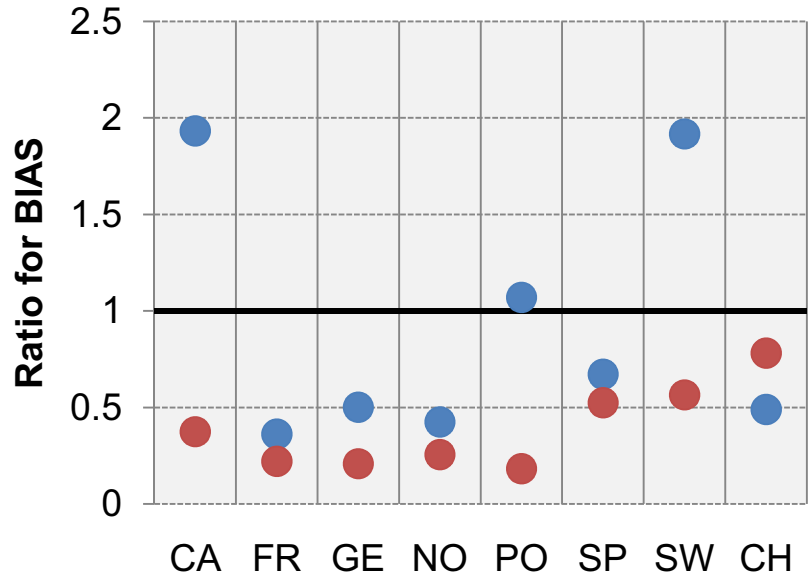


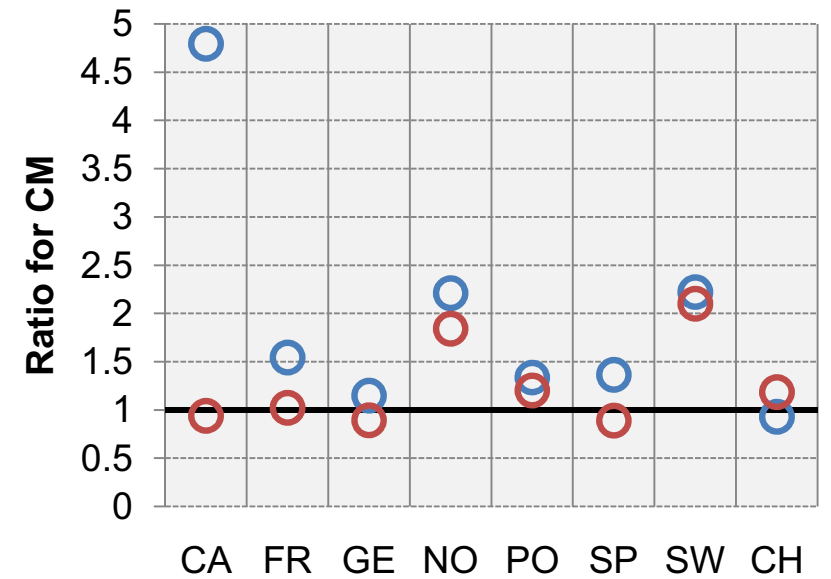
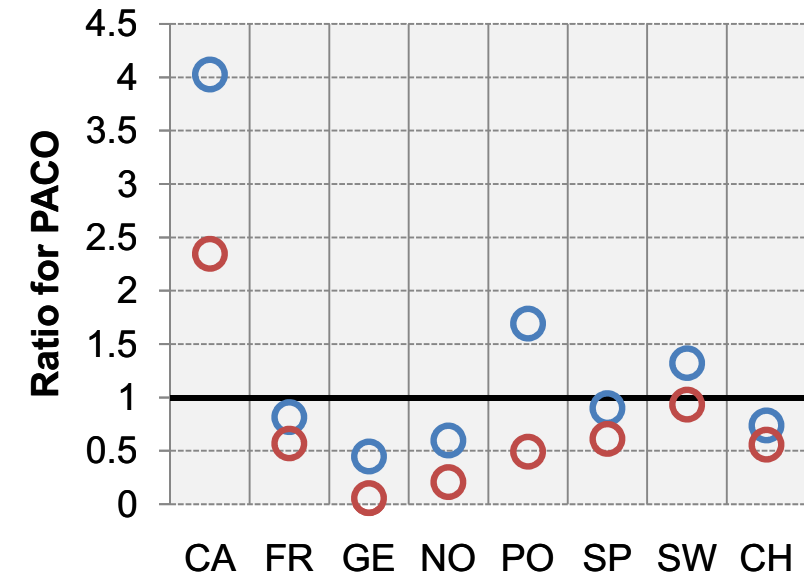
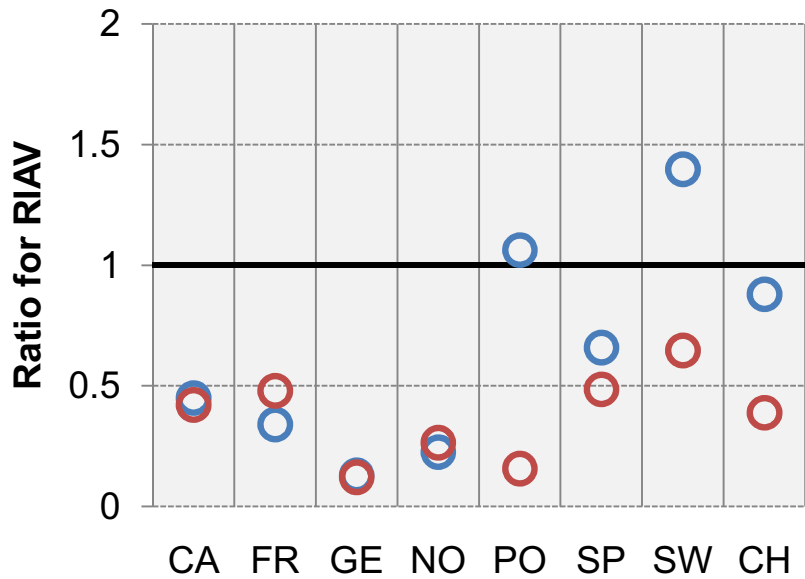
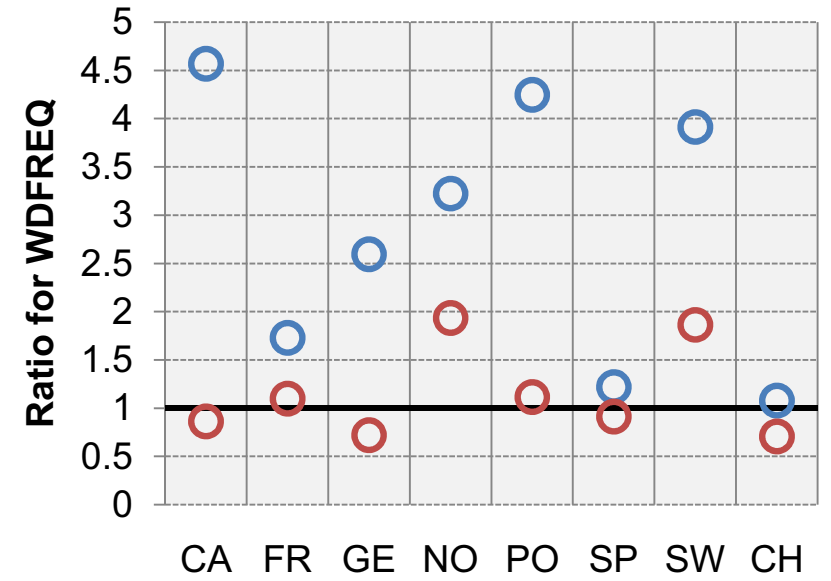
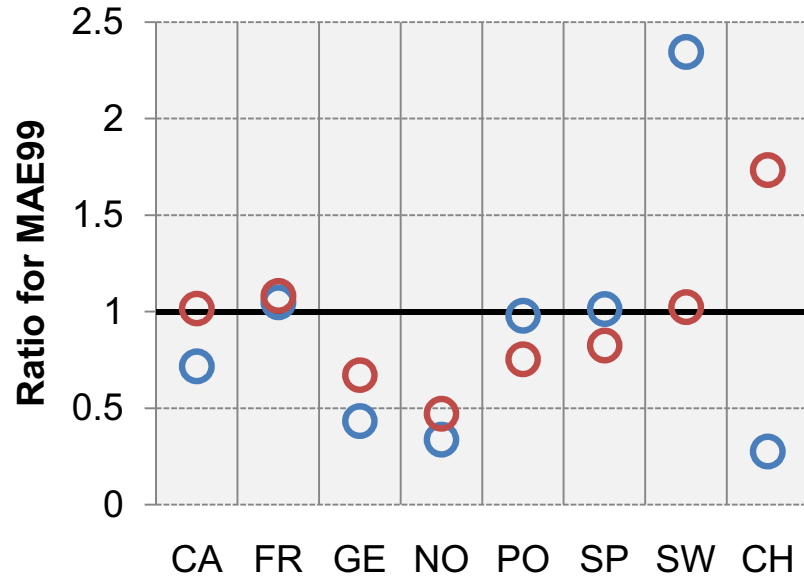
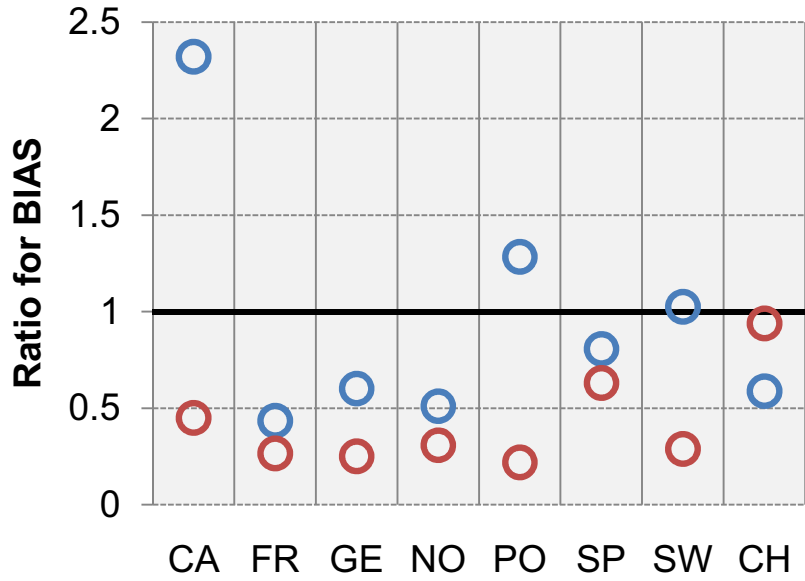
**CM**



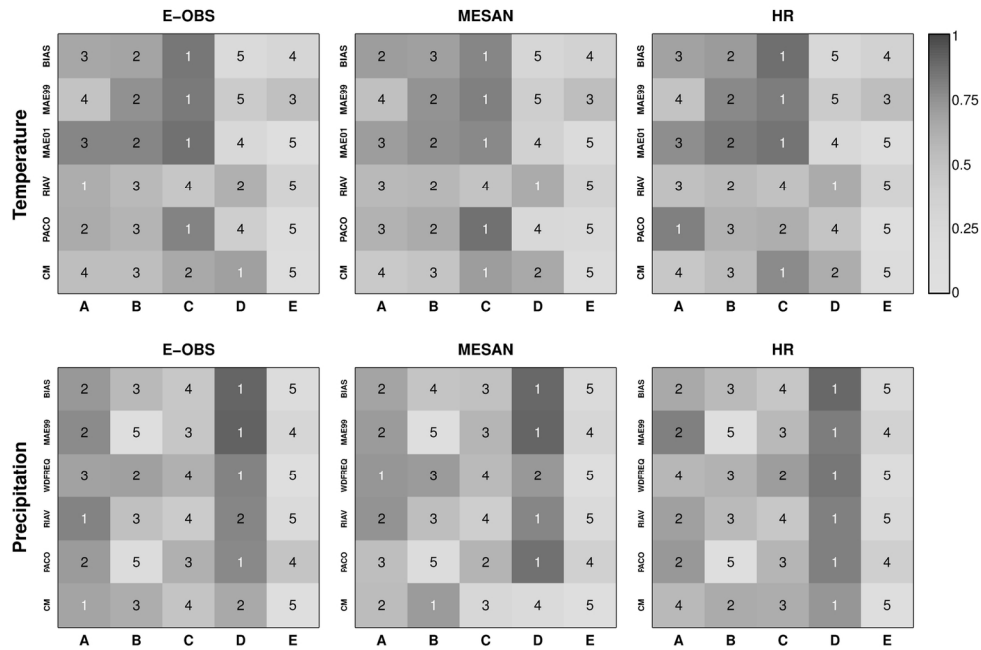
■ E-OBS   
 ■ MESAN   
 ■ HR





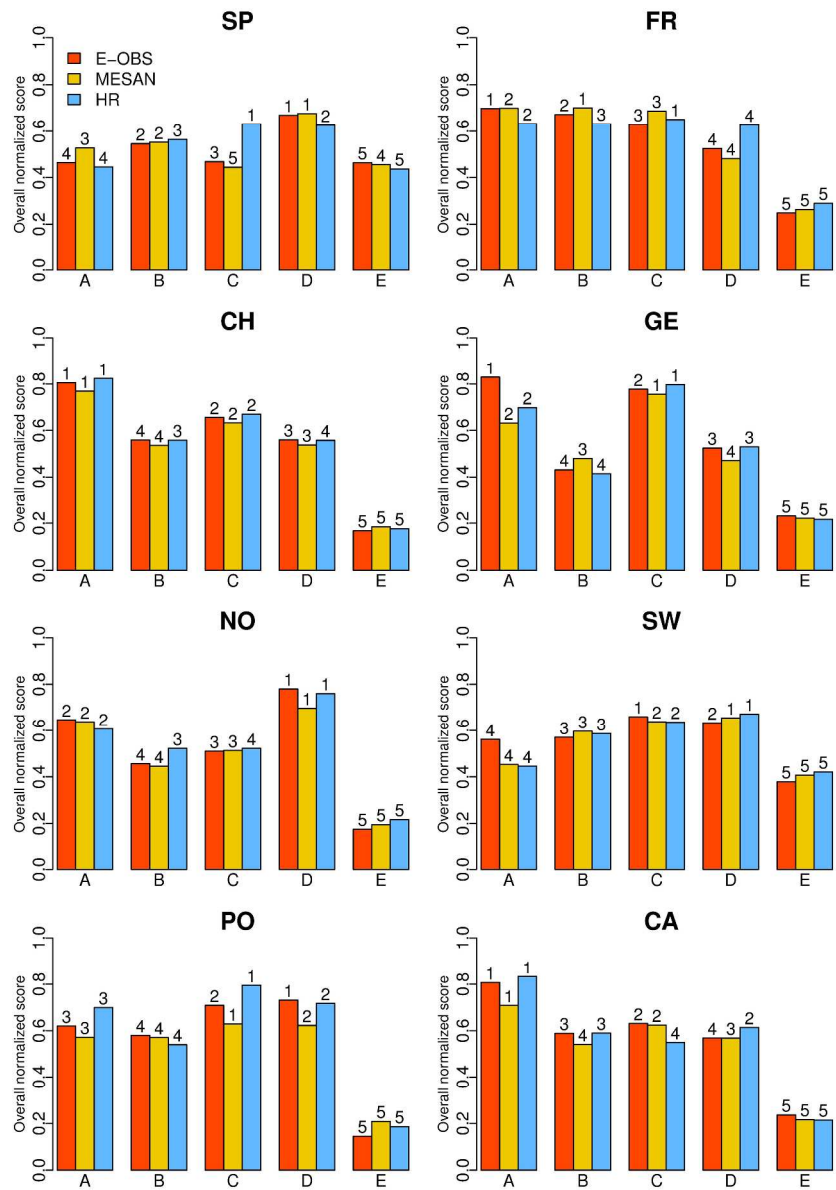






Normalized performance scores (shading) for individual performance metrics, when averaged over all seasons and regions. The upper row shows the results for temperature and the lower row for precipitation. Numbering inside the shaded boxes indicates the actual RCM rank for each case. In each panel, the individual rows indicate the performance metric, the individual columns the five RCMs considered.

152x101mm (300 x 300 DPI)



Overall (combined temperature and precipitation) normalized performance scores for each sub-region. The numbering above the bars indicates the actual RCM ranks separately for each reference dataset.

297x420mm (300 x 300 DPI)

Table 1: Overview on the employed observational reference and RCM datasets. In this work the individual datasets are simply referred to by their abbreviation (last column).

Type of dataset	Details		
Observational reference	Name	Description	Abbreviation
	E-OBS v15	Section 2.1.1	<b>EOBS</b>
	National high-resolution grids	Section 2.1.2	<b>HR</b>
	EURO4M MESAN	Section 2.1.3	<b>MESAN</b>
RCM	Model name and version	Institute/Group	Abbreviation
	CCLM 4.8.17	CLMcom	<b>A</b>
	HIRHAM 5	DMI	<b>B</b>
	WRF 3.3.1F	IPSL-INERIS	<b>C</b>
	RACMO 2.2E	KNMI	<b>D</b>
	RCA 4	SMHI	<b>E</b>