

The VALUE perfect predictor experiment: evaluation of temporal variability

Journal:	<i>International Journal of Climatology</i>
Manuscript ID	JOC-16-0593.R1
Wiley - Manuscript type:	VALUE special issue
Date Submitted by the Author:	17-May-2017
Complete List of Authors:	<p>Maraun, Douglas; University of Graz, Wegener Center for Climate and Global Change Huth, Radan; Charles University, Faculty of Science, Dept. of Physical Geography and Geoecology; Institute of Atmospheric Physics, Dept. of Climatology Gutiérrez, José; National Research Council (CSIC), Instituto de Física de Cantabria; San Martin, Daniel; Predictia Intelligent Data Solutions SL, N.A. Dubrovsky, Martin; Institute of Atmospheric Physics, Dept. of Climatology Fischer, Andreas; Federal Office of Meteorology and Climatology (MeteoSwiss), Climate Services Hertig, Elke; University of Augsburg, Institute for Geography Soares, Pedro; Instituto Dom Luiz, Universidade de Lisboa, DEGGE Bartholy, Judit; Eotvos Lorand Tudomanyegyetem, Department of Meteorology Pongracz, Rita; Eotvos Lorand Tudomanyegyetem, Department of Meteorology Widmann, Martin; University of Birmingham, School of Geography, Earth and Environmental Sciences Casado, María; AEMET, Desarrollo y Aplicaciones Ramos, Petra; Delegacion Territorial de AEMET en Andalucía, Ceuta y Melilla, N.A. Bedia, Joaquin; Predictia Intelligent Data Solutions SL, N.A.</p>
Keywords:	Regional climate, Downscaling, Evaluation, Validation, Temporal variability, Spells, Interannual variability, long-term trends

SCHOLARONE™
Manuscripts

1 The VALUE perfect predictor experiment: evaluation of
2 temporal variability

3 Douglas Maraun¹, Radan Huth^{2,3}, Jose M. Gutierrez⁴, Daniel San Martin⁵,
Martin Dubrovsky³, Andreas Fischer⁶, Elke Hertig⁷, Pedro M. Soares⁸,
Judit Bartholy⁹, Rita Pongracz⁹, Martin Widmann¹⁰, Maria J. Casado¹¹,
Petra Ramos¹² and Joaquin Bedia⁵

¹ Wegener Center for Climate and Global Change, University of Graz,
Brandhofgasse 5, 8010 Graz, Austria

²Dept. of Physical Geography and Geoecology, Faculty of Science,
Charles University; Albertov 6, 128 43 Praha 2, Czech Republic

³ Institute of Atmospheric Physics Czech Academy of Sciences, Bocni II 1401,
141 31 Prague, Czech Republic

⁴ Institute of Physics of Cantabria (IFCA), University of Cantabria,
Avenida de los Castros, Santander 39005, Spain

⁵ Predictia Intelligent Data Solutions SL, Avda. los Castros s/n,
Building I+D S345, 39005, Santander, Spain

⁶ Federal Office of Meteorology and Climatology MeteoSwiss,
Operation Center 1, 8085 Zurich-Airport, Switzerland

⁷ Institute of Geography, Augsburg University, Alter Postweg 118, 86159 Augsburg

⁸ Instituto Dom Luiz, Faculdade de Ciencias, Universidade de Lisboa,
1749-016 Lisbon, Portugal

⁹ Dept. of Meteorology, Eotvos Lorand University, Pazmany st. 1/a,
H-1117 Budapest, Hungary

¹⁰ School of Geography, Earth and Environmental Sciences,
University of Birmingham, Birmingham, B15 2TT, UK

¹¹ Agencia Estatal de Meteorologia (AEMET), C/ Leonardo Prieto Castro, 8
Ciudad Universitaria, 28040 Madrid, Spain

¹² Delegacion Territorial de AEMET en Andaluca, Ceuta y Melilla,
Avda. Americo Vespucio, n 3 bajo., 41092 Sevilla, Spain

4 May 17, 2017

5 **Temporal variability is an important feature of climate, comprising system-**
6 **atic variations such as the annual cycle, as well as residual temporal variations**
7 **such as short-term variations, spells and variability from interannual to long-term**
8 **trends. The EU-COST Action VALUE developed a comprehensive framework to**

9 evaluate downscaling methods. Here we present the evaluation of the perfect pre-
10 dictor experiment for temporal variability. Overall, the behaviour of the different
11 approaches turned out to be as expected from their structure and implementa-
12 tion. The chosen regional climate model adds value to reanalysis data for most
13 considered aspects, for all seasons and for both temperature and precipitation.
14 Bias correction methods do not directly modify temporal variability apart from
15 the annual cycle. However, wet day corrections substantially improve transition
16 probabilities and spell length distributions, whereas interannual variability is in
17 some cases deteriorated by quantile mapping. The performance of perfect prog-
18 nosis statistical downscaling methods varies strongly from aspect to aspect and
19 method to method, and depends strongly on the predictor choice. Unconditional
20 weather generators tend to perform well for the aspects they have been calibrated
21 for, but underrepresent long spells and interannual variability. Long-term tem-
22 perature trends of the driving model are essentially unchanged by bias correction
23 methods. If precipitation trends not well simulated by the driving model, bias
24 correction further deteriorates these trends. The performance of PP methods to
25 simulate trends depends strongly on the chosen predictors.

26 1 Introduction

27 Downscaling is a common - often necessary - step in assessing regional climate change and
28 its impacts: the resolution of global coupled atmosphere-ocean general circulation models
29 (GCMs) is typically too coarse to represent many regional- or local-scale climate phenomena.
30 Therefore the output of GCMs is downscaled to provide high resolution simulations over a
31 limited target area. The EU Cooperation in Science and Technology (COST) Action ES1102
32 VALUE was established to comprehensively evaluate different downscaling methods (Maraun
33 et al., 2015). Three experiments have been defined: a so-called perfect predictor experiment
34 to isolate downscaling skill in present climate; a GCM predictor experiment to evaluate the
35 overall skill to simulate present-day regional climate; and a pseudo reality experiment to
36 evaluate the skill of downscaling methods to represent future climates.

37 In a community effort, researchers from 16 European institutions participated in the per-
38 fect predictor experiment, and more than 50 different statistical downscaling methods have
39 been evaluated at 86 stations across Europe. The evaluation comprises the representation of
40 marginal aspects (such as the mean or variance; (Gutiérrez and coauthors, 2017)), temporal
41 aspects (such as spell length distributions; this contribution), spatial aspects (such as spatial
42 decorrelation lengths; (Widmann and coauthors, 2017)), and multivariable aspects (such as
43 the relationship between temperature and precipitation; Page et al., in preparation). Extreme
44 events as well as an evaluation conditional on relevant synoptic and regional phenomena have
45 been, owing to their importance, considered separately by Hertig and coauthors (2016) and
46 Soares and coauthors (2017). Here we present the evaluation of temporal aspects.

47 To illustrate different aspects of temporal variability, Figure 1 shows a selected year of
48 precipitation at the participating rain gauge in Graz, Austria. On 18th of July (orange spike),
49 several districts were flooded. The city's streams burst their banks following the heavy rainfalls
50 prior to the event, but a major contributor was the long wet spell in the end of June (red

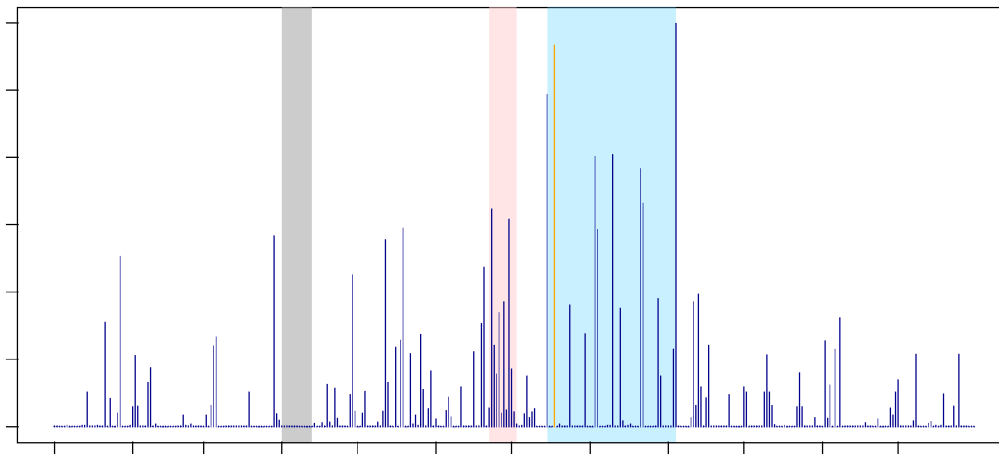


Figure 1: Daily precipitation totals in Graz, 2009. Shading: see text.

51 shading). Southeast of Graz, the overall event caused several thousand landslides. Total
52 rainfall in June exceeded the climatological mean by more than 60%. Also annual rainfall was
53 about 47% higher than normal (Klein Tank et al., 2002), indicating substantial interannual
54 variability. A pronounced seasonality of all aspects of precipitation is directly apparent. In
55 late winter and early spring, precipitation amounts are low compared to summer. Also the
56 probability of consecutive wet days is low resulting in long dry spells (grey shading). Most
57 dry-wet and wet-wet transitions occur in late spring and early summer, the highest rainfall
58 amounts are observed in late summer (blue shading).

59 In general, temporal variability involves a wide range of time scales, from the diurnal cycle
60 through day-to-day variations, spells (dry, wet, warm, cold, etc.), and interannual variations
61 to long-term trends. The variability can be broadly separated into systematic variations
62 - the diurnal and annual cycle as well as forced long-term trends - and residual temporal
63 variations, whose characteristics are determined by the large-scale driving processes and by
64 local memory. For instance, temporal dependence in precipitation may stem directly from
65 memory caused by soil-moisture feedbacks, or indirectly from the duration of passing cyclones
66 and anti-cyclones. Temporal aspects of local climate are often essential for impact studies in
67 various sectors such as water (e.g., preconditions of flooding, Froidevaux et al. (2015); dry
68 spells, Stoll et al. (2011)), agriculture (e.g., dry spells Calanca (2007); seasonality, Rosenzweig
69 et al. (2001)), health (Semenza et al., 1996, e.g., heatwaves) and energy (Rosenzweig et al.,
70 2011, e.g., seasonality).

71 In VALUE we evaluate the performance of different downscaling methods to represent
72 temporal variability. Apart from dynamical downscaling with regional climate models (RCMs,
73 Rummukainen, 2010), different statistical approaches exist (Fowler et al., 2007; Maraun et al.,
74 2010; Wilks, 2010; Maraun, 2016): perfect prognosis (PP) statistical downscaling methods,
75 which are calibrated purely on observations and typically take their predictors from large-scale
76 fields of the free atmosphere; model output statistics (MOS) methods, which are calibrated
77 between model data and observations (in climate science, these are typically bias correction

78 methods); and unconditional weather generators, which are calibrated on local data and do
79 not include any meteorological predictors.

80 The basic driver of the residual, regional-scale temporal variability is the propagation of
81 planetary and synoptic waves, which is essentially prescribed by GCMs. This continental-scale
82 variability is modulated by regional-scale dynamical processes, influences of the orography, and
83 feedback mechanisms such as soil-moisture-temperature, soil-moisture-precipitation feedbacks
84 and snow-albedo feedbacks (Schär et al., 1999; Seneviratne et al., 2006; Fischer et al., 2007;
85 Hall et al., 2008). As a result, regional-scale temporal variability simulated by RCMs may
86 diverge from the prescribed large-scale variability (Alexandru et al., 2007). Local temporal
87 variability is often - in particular for precipitation and wind - not fully determined by larger-
88 scale variability, but exhibits additional - essentially random - fluctuations. PP statistical
89 downscaling inherits the variability of the large-scale predictors and typically does not add
90 any local short-term variations. Some methods, however, explicitly model local variability by
91 randomisation (von Storch, 1999; Chandler and Wheeler, 2002; Volosciuk et al., 2017). Such
92 stochastic models might simply generate white noise, but may also include weather genera-
93 tors (see below) to model short-term temporal dependence by Markov-chain-type components
94 (Maraun et al., 2010). Also bias correction typically does not explicitly add local temporal
95 variability to the driving model, but only subtly modulates temporal variability via its effect
96 on the marginal distribution. For instance wet day frequencies are adjusted, which indirectly
97 affects the representation of spells (Rajczak et al., 2016). Some bias correction methods also
98 attempt to explicitly adjust the temporal structure (Vrac and Friederichs, 2015; Cannon, 2016,
99 e.g.) but at the cost of destroying the temporal consistency with the driving dynamical model.
100 Unconditional weather generators (i.e., weather generators that do not use meteorological pre-
101 dictors) do not provide sequences which are synchronised with the driving models. Instead,
102 the only temporal structure they represent is explicitly modelled, typically by Markov chains
103 (Maraun et al., 2010). Most statistical models - PP and MOS - have an explicit description of
104 the annual cycle, e.g., by being calibrated to each calendar day, month or season individually,
105 or (in case of PP) by including the day-of-the year as predictor.

106 Of the temporal aspects studied in this paper, perhaps the annual cycle has been the
107 most frequent target of validation: many RCM studies as well as studies of both kinds of
108 statistical downscaling (PP and MOS) and of WGs include a validation of the annual cycle,
109 although it usually is not their main topic (e.g. Frei et al., 2003; Moberg and Jones, 2004;
110 Kilsby et al., 2007; Turco et al., 2011; Schindler et al., 2007; Soares et al., 2012; Warrach-
111 Sagi et al., 2013; Kalognomou et al., 2013; Martynov et al., 2013; Keller et al., 2015; Favre
112 et al., 2016). Also studies evaluating precipitation (dry/wet) spells and precipitation transi-
113 tion probabilities (wet/wet, dry/wet) as well as interannual variability have been relatively
114 numerous (e.g. Semenov et al., 1998; Charles et al., 1999; Giorgi et al., 2004; Kilsby et al.,
115 2007; Jacob et al., 2007; Schmidli et al., 2007; Frost et al., 2011; Bürger et al., 2012; Turco
116 et al., 2011; Hu et al., 2013; Gutmann et al., 2014; Keller et al., 2015; Rajczak et al., 2016).
117 Much less attention has, on the other hand, been paid to validation of temperature spells and
118 day-to-day temperature changes; only a few studies have been published that focus on these
119 characteristics (Huth et al., 2001; Bürger et al., 2012; Vautard et al., 2013; Huth et al., 2015;
120 Lhotka and Kyselý, 2015).

121 The vast majority of validation studies addressing also temporal issues focused on a single
122 downscaling approach or, at best, provide a comparison for models from one family such as
123 Kotlarski et al. (2014); Gutmann et al. (2014). Exceptions are Wilby et al. (1998), who where

124 the first to systematically evaluate temporal aspects in PP methods and unconditional weather
125 generators; the STARDEX project, which assessed temporal aspects of extreme events in PP
126 and a simple MOS method (Haylock et al., 2006; Goodess et al., 2010); the study by Frost
127 et al. (2011), who compared the representation of spell lengths and interannual variability in
128 an RCM, a bias correction method, a PP method and two weather generators; the study by
129 Hu et al. (2013), who carried out a similar intercomparison for a PP method and two weather
130 generators; the study by Bürger et al. (2012), who compared extreme spells in several PP
131 and MOS methods; and the recent study by Huth et al. (2015), which investigated temporal
132 aspects in both statistical and dynamical downscaling methods. But all these studies still
133 include only a rather limited range of methods.

134 Even though extremely important for climate change studies (Pielke and Wilby, 2012),
135 evaluation studies of trends in downscaled data are scarce (Benestad and Haugen, 2007; Lorenz
136 and Jacob, 2010; Bukovsky, 2012; Ceppi et al., 2012; Huth et al., 2015). These studies broadly
137 indicate a rather limited ability of downscaling methods to reproduce trends.

138 In brief, a substantial research gap exists. The performance of many downscaling and
139 bias correction methods to represent temporal aspects - both individually and relative to
140 each other - is largely unknown. This study takes a first step to close this gap. In a perfect
141 predictor experiment we analysed the performance of one raw RCM and 48 statistical methods
142 to represent day-to-day variability, spells, seasonality, interannual and long-term variability
143 including trends. Aspects of temporal variability specifically addressing extreme events, such
144 as long heatwaves or meteorological drought, are addressed in the companion paper on extreme
145 events (Hertig and coauthors, 2016, in this issue). The considered experiment was conducted
146 for daily values, hence we cannot evaluate sub-daily variations.

147 VALUE is a community effort, the participation in this experiment (and its evaluation) was
148 unpaid. The participating methods thus form an ensemble of opportunity. In particular no
149 systematic set of predictor variables or domains has been prescribed. Thus statements about
150 optimal predictor choice are limited to a few comparisons of similar (or identical) methods
151 with different predictors. A detailed set of metadata has, however, been collected for all
152 participating methods. These meta data describe structural aspects of all methods and often
153 allow for quite detailed interpretations of the individual performance. In the paper we will
154 discuss selected examples in more detail, and additionally give a broad overview of the different
155 model families. The metadata and complete results for individual methods are available from
156 the VALUE portal www.value-cost.eu/validationportal for further investigation.

157 The aim of the perfect predictor experiment is to evaluate the isolated skill of the raw RCM
158 and the statistical models. Consequently, this study cannot give a conclusive assessment of
159 the skill to simulate regional future climates. The skill of a full regional modelling system,
160 comprising the full modelling chain from GCM to RCM and/or statistical model, as well as
161 the downscaling performance in future climates will be considered in additional experiments
162 (Maraun et al., 2015).

163 In the following section we will briefly review the experimental setup, the considered di-
164 agnostics and the participating methods. In Section 3 we will present the results for different
165 diagnostics and methods. An overall discussion of the results will follow in the final section.

166 2 Experiment, Diagnostics and Methods

167 The experimental design follows the VALUE perfect predictor experiment with station data
168 as target. As (approximately) perfect predictors and perfect boundary conditions, we use
169 ERA-Interim data from 1 Jan 1979 to 31 Dec 2008 (Dee et al., 2011). The MOS methods
170 use ERA-Interim data at their native resolution of 0.75° as input, the PP methods ERA-
171 Interim predictors at 2° , which resembles a typical GCM resolution. Furthermore, most MOS
172 methods also use ERA-Interim, downscaled with the RCM RACMO (van Meijgaard et al.,
173 2008), as input to represent a typical RCM bias correction situation. Apart from the resolution,
174 some important differences between these two MOS settings exist: in the first case, internal
175 variability at the grid-box scale is closely tied to real world internal variability, whereas the
176 RCM develops its own internal variability within the RCM domain. Furthermore, observed
177 temperatures have been assimilated into the ERA-Interim reanalysis; the resulting predictors
178 are thus essentially bias free at the grid-box scale and differences with station observations
179 mainly result from the scale gap. RCM temperatures inside the domain, however, are only
180 mildly constrained by the boundaries and are thus typically affected by biases. Precipitation
181 is in both cases calculated by model parameterisations, without any reference to observed
182 precipitation. It is thus affected by scale-gap and biases.

183 As predictand data, time series from 86 stations from the publicly available ECA data base
184 were used (Klein Tank et al., 2002). These stations were selected to cover the different Euro-
185 pean climates, covering mediterranean, maritime, continental, alpine and sub-polar climates.
186 For details refer to Gutiérrez and coauthors (2017) and the supplementary information.

187 In this manuscript, we consider daily maximum and minimum temperature and daily
188 precipitation only. A dedicated analysis of other variables will be carried out separately for
189 a set of stations in Germany (Page et al., in preparation). For the statistical methods a five-
190 fold cross validation with non-overlapping 6-year blocks is carried out. Further details about
191 the protocol can be found in Maraun et al. (2015), Gutiérrez and coauthors (2017) and on
192 www.value-cost.eu/validation#Experiment_1a.

Index	Variables	Performance measure	Resolution	Description
short-term variability				
ACF1	T_{max} , T_{min}	bias	seasonal	lag-1 autocorrelation
ACF2	T_{max} , T_{min}	bias	seasonal	lag-2 autocorrelation
WWprob	precipitation	bias	seasonal	probability of wet-wet transition
WDprob	precipitation	bias	seasonal	probability of wet-dry transition
Spells				
WarmSpellMean	T_{max}	bias	seasonal	mean of the warm (> 90th percentile) spell length distribution
ColdSpellMean	T_{min}	bias	seasonal	mean of the cold (< 10th percentile) spell length distribution
WetSpellMean	precipitation	bias	seasonal	mean of the wet (≥ 1 mm) spell length distribution
DrySpellMean	precipitation	bias	seasonal	mean of the dry (< 1mm) spell length distribution
Interannual to long-term variability				
VarY	T_{max} , T_{min} , precipitation	rel. error	seasonal	variance of seasonally/annually averaged data
Cor.1Y	T_{max} , T_{min} , precipitation	bias	seasonal	correlation with observations of seasonally/annually averaged data
Cor.7Y	T_{max} , T_{min} , precipitation	correlation	seasonal	correlation with observations of seasonally/annually averaged and filtered data
Trend	T_{max} , T_{min} , precipitation	trends themselves	seasonal	long-term (relative) trend of seasonally/annually averaged data
Annual cycle				
AnnualCycleAmp	T_{max} , T_{min}	bias	annual	Amplitude of the annual cycle
AnnualCycleRelAmp	precipitation	rel. error	annual	Relative amplitude of the annual cycle
AnnualCyclePhase	T_{max} , T_{min}	circular bias	annual	Phase of highest peak ²

Table 1: Diagnostics considered. Diagnostics only shown in the supplementary information are plotted in grey. For details see <http://www.value-cost.eu/validationportal/app#!indices> and click on “details” for the underlying R-Code (note that registration is required).

193 Table 1 lists the diagnostics we considered: the indices to measure a specific aspect of
194 temporal variability, the corresponding performance measure to quantify the mismatch with
195 observations and the temporal resolution (seasonal, annual) at which the evaluation has been
196 carried out. In two cases, we assessed correlations between observed and downscaled local
197 time series, namely at the interannual and seven year time scales. In this case, the diagnostic
198 consists of a performance measure - the correlation - only.

199 Detailed descriptions of these diagnostics can be found in the supplementary information.
200 The code used to calculate these diagnostics is available from
201 <http://www.value-cost.eu/validationportal/app#!indices> (registration required).

202 In this analysis, we compare methods from the PP, MOS and unconditional weather gen-
203 erator approaches with raw ERA-Interim output, and dynamically downscaled ERA-Interim.
204 Tables 2 and 3 list the methods participating in the experiment (many methods are iden-
205 tical for the different variables, but in several cases differences exist in the implementation
206 for different variables. Therefore, we decided not to list the methods in a single table). The
207 MOS methods are listed prior to the PP methods to ease comparison with the raw RCM and
208 ERA-Interim data.

209 PP methods are calibrated purely on observed predictors and predictands. The statistical
210 model is then applied to climate model predictors. In a climate change context, the approach
211 is based on three major assumptions Maraun and Widmann (2018): first, that the GCM
212 predictors are perfectly simulated (hence the name) in present and future climate. As a
213 consequence, predictors are typically taken from large-scale fields of the free atmosphere.
214 Second, the predictors should be informative of local variability and climate change. And
215 third, the model structure should well describe local variability, and allow for at least moderate
216 extrapolations under climate change. Our evaluation experiment employs perfect predictors
217 to isolate downscaling skill in present climate. It can therefore be used to assess whether the
218 chosen predictors are informative of local variability and observed changes, and whether the
219 model structure well describes observed local variability and changes. The perfect prognosis
220 assumption and performance under future climate change, however, cannot be assessed.

221 The participating PP methods broadly represent widely used approaches - analogue, re-
222 gression and weather-type methods. Some of regression methods apply variance inflation
223 (MLR-ASI, MLR-AAI, GLM-P), some are stochastic (see Tables). The ESD methods down-
224 scale at the monthly scale, thus no diagnostics are considered that involve daily values. The
225 ESD-EOF implementation differs from the standard ESD version in that the predictand values
226 are filtered by PCA Benestad et al. (2015b).

227 All stochastic methods use, conditionally on the predictors, independent noise, i.e., they
228 do not have an explicit Markov component implemented to simulate short-term persistence.
229 For precipitation, some of the participating PP methods have been included for illustrative
230 purposes only (MLR-RAN, MLR-RSN, MLR-ASW, MLR-ASI). In fact, it is well known that
231 simple multiple linear regression methods are not suitable to model daily precipitation. Yet
232 they do participate in the intercomparison to highlight the problems associated with them
233 (marked in grey in Table 3). Two of the stochastic methods (GLM and SWG) are based
234 on generalised linear models, with a logistic regression for the occurrence process, and a
235 generalised linear regression on the gamma distribution parameters for the amounts process.
236 GLM-WT and WT-WG condition the distribution parameters for occurrence and amounts on
237 weather types.

238 MOS methods are calibrated between model simulations and observations. The approach

239 can thus in principle adjust biases (in fact, in climate science, these are almost exclusively bias
240 correction methods, i.e., predictor and predictand have the same physical dimension), but has
241 to be calibrated individually to the chosen model. MOS is based on three major assumptions
242 (which make up the so-called stationarity assumption), similar to those of the PP approach
243 Maraun and Widmann (2018) : first, the predictors have to be credibly (but not necessarily
244 bias free) simulated. Second, the predictors need to be representative of the local variable.
245 And third, as in PP, the structure of the transfer function needs to be suitable. Again, the
246 first assumption cannot be tested with perfect predictors, only the second and third, and only
247 for present day climate.

248 The participating MOS methods comprehensively span the range of widely used methods,
249 and also cover some more experimental recent developments such as stochastic bias correction
250 (VGLMGAMMA Wong et al., 2014). None of the participating MOS methods modifies resid-
251 ual temporal dependence directly, but only indirectly via changes in the marginal distribution.
252 The CDFt method calibrates a statistical distribution also in the validation period. As this
253 is only 6 years in our experiment (in a climate change experiment, one would typically use a
254 30 year time slice), we expect a broad spread for the resulting performance measures due to
255 sampling variability.

256 Unconditional weather generators are not conditioned on meteorological predictors, but
257 stochastically simulate marginal and temporal aspects, sometimes also spatial. They are
258 calibrated to observed weather statistics. Under climate change, the model parameters (or
259 the observed weather statistics) are adjusted by so-called change factors derived from climate
260 models. The underlying assumptions are thus similar to those for MOS Maraun and Widmann
261 (2018): first, the change factors have to be credibly simulated, and all relevant change factors
262 have to be included; second, the simulated change factors have to be representative of local
263 changes; and third, the model structure has to be suitable. In the chosen experiment, no
264 change factors are applied between calibration and validation period; thus only the suitability
265 of the model structure can be evaluated. Some climatic statistics may have changed between
266 calibration and validation period, but resulting systematic biases cancel out under cross-
267 validation.

268 The SS-WG and MARFI unconditional weather generators are of the Richardson type
269 Richardson (1981), i.e., they use a Markov chain to simulate precipitation occurrence, and
270 an autoregressive model to simulate temperature. A major difference between the two is
271 the wet-day threshold: the SS-WG uses 1 mm, the MARFI models use 0.5 mm (note that
272 the evaluation indices are in any case based on a 1 mm threshold). The GOMEZ weather
273 generators are based on resampling.

274 Diagnostics have been calculated for each method and each station. They can be down-
275 loaded from the VALUE portal (www.value-cost.eu/validationportal/app#!validation).
276 For stochastic methods, an ensemble of 100 realisations have been uploaded. The performance
277 measures have been derived for each realisation and then averaged across the ensemble.

278 When interpreting the evaluation results, it has to be acknowledged whether a specific
279 index is calibrated or emerges from the model. For instance, a good representation of the
280 annual cycle could result from including meteorological predictors that describe the annual
281 cycle, or trivially from fitting a statistical model separately to each month. In particular,
282 weather generators by construction resemble many marginal and temporal aspects. In this
283 study, only spell lengths and interannual variability are not calibrated. In Tables 2 and 3 we
284 therefore also list whether short-term dependence (AC) and seasonality (SE) are calibrated or

285 not. For further details on the contributing methods see Gutiérrez and coauthors (2017) or
286 the VALUE portal (www.value-cost.eu/validationportal/app#!downscalingmethod).

287 3 Results

288 Figure 2 illustrates selected temporal aspects for precipitation in Graz, Austria, and how
289 corresponding model performance has been quantified in this study. The top panel shows the
290 dry spell length distribution. Observations are shown in bold solid black, the results for five
291 different statistical methods are shown in color. Methods in red and orange are MOS, in blue
292 PP, and the method shown in magenta is an unconditional weather generator. One index that
293 can be derived from the distribution is the mean spell length (which is quantified in this study
294 for all the participating methods and all selected weather stations). Dashed vertical lines show
295 this index for observations and statistical models. The performance of a model is given by the
296 difference between the modelled and observed mean, i.e., the mean spell length bias. Similarly,
297 the bottom panel shows the annual cycle of daily mean precipitation. Here, two indices are
298 considered: first, the relative amplitude (for temperature the absolute amplitude) defined as
299 the difference between maximum and minimum value (horizontal dashed lines), relative to the
300 mean of these two values. Second, the phase of the annual cycle, defined as the day of the
301 annual cycle maximum⁴ (vertical dashed lines). The performance for the first is measured as
302 the relative error between modelled and observed relative amplitude, for the second as the
303 circular bias between modelled and observed phase (circular in the sense that the difference
304 between, say, 31st of December and 1st of January is -1 day, not 364 days).

305 In the following, we present the results, separately for temperature and precipitation. To
306 keep the number of figures at a reasonable level, we selected a suite of relevant diagnostics for
307 short-term variability, spells, monthly to interannual variability, and the annual cycle. Often,
308 only one season is shown, in case of temperature, only either daily minimum or maximum
309 temperature. A more comprehensive catalog of plots can be found in the supplementary
310 information. The figures for all diagnostics are organised similarly, see Fig. 3 as an example.
311 In this example, one diagnostic is shown for daily maximum and minimum temperature. In
312 the top row, the observed indices are shown - here auto-correlation of daily maximum (left)
313 and minimum (right) temperatures. Note that correlations on interannual and 7-year time
314 scales have no corresponding observed indices, consequently no maps are drawn. The two
315 panels below show the performance measures for these indices (top: maximum temperature,
316 bottom: minimum temperature). Each box-whisker-plot represents one method: the raw
317 driving data (ERA-Interim at the 2° resolution used as predictor for PP methods, at the
318 native 0.75° resolution and the RACMO2 RCM), the MOS methods, the PP methods and the
319 unconditional weather generators. The individual box-whisker-plots summarise the results for
320 all 86 stations: the boxes give the 25%-75% range, the whiskers the maximum value within
321 1.5 times the interquartile range; values outside that range are plotted individually. The thick
322 colored horizontal bars show the medians for the individual PRUDENCE regions (Christensen
323 and Christensen, 2007). Note that the number of stations entering these calculations differs
324 from region to region (ranging from 3 in France to 21 in Scandinavia, typically around 10).
325 A red asterisk indicates that values lie outside the plotted range. Results for individual

⁴In some cases, the annual cycle of precipitation has two maxima. We will discuss below how the phase is defined in this case.

326 stations are - depending on the index - substantially affected by noise, but the median over all
327 considered stations in general provides a robust estimate of the overall performance of a given
328 method. Furthermore, the diagnostic is solely defined between observations and simulations,
329 thus no observed indices exist.

330 For a given index, all methods are shown for which the index may sensibly be calculated.
331 That is, methods producing only monthly output are not shown for any indices based on daily
332 values. Otherwise, all indices are presented, even though a method might not be designed to
333 reproduce them. Such results are not intended to denounce specific methods, but rather to
334 highlight the consequences of using a method in such a context. These situations will be made
335 explicit to avoid misinterpretation of the results.

336 As mentioned in the introduction, the methods participating in the experiment form an
337 ensemble of opportunity. Also we have a list of candidate predictors for each method, but
338 the actually selected set of predictors might be much lower for individual stations. To fully
339 attribute differences in model performance to the approach, the particular implementation
340 and the choice of predictors, dedicated sensitivity studies would be required. In many cases,
341 conclusions may be drawn for groups of methods. For instance, all analog methods often be-
342 have similarly independent of the different predictors and implementations. Thus, conclusions
343 about analog-type methods as a whole can often be drawn. A discussion of differences within
344 this type, however, would be very speculative, because the individual methods often differ
345 both in the implementation and choice of predictors. The level of detail in our interpretation
346 will thus differ from case to case. In some cases, any discussion would be too speculative - we
347 then restrict ourselves to a description of the findings.

348 3.1 Temperature

349 **short-term variability** Figure 3 shows the results for lag-1 autocorrelation of summer
350 daily maximum and minimum temperature as a measure of short-term persistence. The top
351 row shows observations for daily maximum (left) and minimum (right) temperature. The
352 corresponding plots for winter can be found in the supplementary information. For T_{max} ,
353 summer persistence is relatively evenly distributed across Europe; for T_{min} , persistence is
354 notably lower over many regions. The bottom panels show the performance of the individual
355 models.

356 The spatial averaging of ERA-Interim results in a moderate overestimation of summer
357 persistence of T_{max} (upper panel), these biases are reduced by the RCM. Almost all MOS
358 methods inherit the skill of the predictor data set, in particular the added value of the RCM.
359 The regression based MOS method (MOS-REG) includes averaging across several grid boxes
360 and thus overestimates persistence. All analog methods underestimate persistence of temper-
361 ature. The reason might be twofold: first, the spatial predictor variability might be strongest
362 for circulation-based predictors. Thus, analogs may be selected that best constrain circula-
363 tion (and in turn precipitation, see Section 3.2). And second, large-scale analogs might be
364 sufficiently dissimilar at local scales to deteriorate day-to-day variations. Understanding this
365 problem requires further detailed analysis. The ANALOG-ANOM method uses predictors
366 defined at a continental scale, which likely explains the low performance.

367 As expected, all deterministic regression models overestimate persistence, as not all local
368 variability is explained by large-scale predictors. This problem cannot be mitigated by inflated
369 regression (MLR-ASI, MLR-AAI). All stochastic regression models randomise with white noise

370 (MLR-ASW, MLR-AAW; though conditional on the predictors) and thus underestimate per-
371 sistence. The low performance of the SWG method may partly be explained by the use of
372 continental-scale predictors in combination with a stochastic white-noise randomisation. The
373 WT-WG method performs worst, as it is stochastic and additionally uses only sea level pres-
374 sure as predictor. For the Iberian Peninsula and the UK, ERA-Interim overestimates summer
375 persistence of T_{max} , the RCM reduces the bias. Conversely, for Eastern Europe ERA-Interim
376 is almost bias free, but the RCM reduces persistence. This performance is again inherited by
377 many statistical methods.

378 For T_{min} (lower panel), the performance is consistently worse for all approaches, with a
379 strong tendency to overestimate summer persistence. The RCM, however, performs slightly
380 worse than ERA-Interim. The relative performance across most other methods is similar to
381 that for T_{max} . The ISIMIP method, driven with ERA-Interim, is a notable exception - it
382 has the lowest bias of all MOS methods. Most MOS methods leave the persistence bias es-
383 sentially unchanged, the methods driven with reanalysis data have a lower bias, the methods
384 driven with the RCM a higher. Interestingly, however, some QM-based bias correction meth-
385 ods moderately improve the representation of persistence indirectly by adjusting marginal
386 distributions. The persistence of summer T_{min} is overestimated in the British Isles. But in
387 contrast to the overall behaviour, this bias is reduced by the RCM (and again, this reduction
388 is inherited by the MOS methods). The performance for most methods is best in the Alps.

389 **Spells** Overall, the performance to simulate spells is similar to the performance to simulate
390 short-term variability. The results for summer temperature spells are shown in Figure 4,
391 measured in terms of the mean spell length. Recall that temperature-related spells are not
392 defined by exceedances of absolute thresholds (e.g., 30°C), but by the 90th percentile of
393 daily maximum temperature, which varies from station to station and will be much lower in
394 Scandinavia than in the Mediterranean (Table 1). The longest summer warm spells occur
395 in Scandinavia, the shortest in the western Mediterranean. Summer cold spells are generally
396 much shorter shortest in Northern Europe, and longest in the Mediterranean.

397 ERA-Interim simulates slightly too long warm spells of T_{max} (upper panel), in particular
398 for the area averaged version. The RCM, again, adds value. MOS inherits the predictor
399 performance (by construction, as the percentile-based spells are invariant to bias correction).
400 Owing to the predictor averaging, the regression based MOS (MOS-REG) again performs
401 considerably worse. Also the behavior of the PP methods is broadly consistent with that
402 for short-term persistence: analog methods and stochastic white noise methods (MLR-ASW,
403 MLR-AAW, WT-WG, SWG) simulate too short spells. This holds in particular WT-WG,
404 driven only with sea level pressure. Weather generators slightly underestimate mean spell
405 lengths, in particular those who underestimate short-term persistence. Persistence of summer
406 warm spells of T_{max} is consistently overestimated over the Mediterranean, a bias which is
407 much improved by the RCM.

408 The persistence for summer cold spells of T_{min} (lower panel), consistent with the results
409 for short-term persistence, is generally too high. The RCM deteriorates the performance of
410 ERA-Interim. This performance is, again trivially, unchanged by the MOS methods. The
411 PP methods perform similar as for warm spells, though with a tendency towards higher
412 persistence. All weather generators perform well, consistent with the results for short-term
413 persistence. Cold spells of summer T_{min} are too long for the British Isles and (but to a lesser
414 extent) the Mediterranean. Performance is best for the Alps.

415 **Seasonality** The amplitude of the annual cycle of T_{max} (Figure 5) is small towards the
416 Atlantic and the Mediterranean, and large in the continental climates of eastern Scandinavia
417 and Eastern Europe. It peaks in July in continental central and eastern Europe, and slightly
418 later in August towards the Atlantic. ERA-Interim slightly underestimates the amplitude
419 of the seasonal cycle (upper panel) - likely linked to its resolution, as the further averaging
420 increases the bias. The RCM in general adds value, but also increase spread across stations.
421 Being seasonally trained, most MOS methods trivially capture the annual cycle well. Note,
422 however, that also the quantile mapping methods without an explicit annual cycle perform
423 well (GPQM, EQM, EQM-WT) for most stations. The authors do not understand the strong
424 drop in performance of the MOS-REG method when driven with the RCM instead of ERA-
425 Interim. Most PP methods perform reasonably well, even those without seasonal training,
426 because the physical link between the predictors (including temperature) and the predictand
427 is close. Only the WT-WG method sticks out: it is not seasonally trained and uses only
428 sea level pressure as predictor. Thus, seasonality in circulation patterns is captured, but not
429 the changes in temperature within these patterns. The weather generators perform well by
430 construction.

431 The phase of the seasonal cycle (lower panel) is captured by most methods. ERA-Interim
432 peaks a day too late, the RCM increases the spread across stations. MOS methods perform
433 well, even those with an explicit model of the seasonal cycle (GPQM, EQM, EQM-WT) are
434 within ± 2 days (apart from the MOS-REG method, when driven with the RCM). The analog
435 methods perform reasonably well, although the version without seasonal training (ANALOG)
436 has a comparably broad spread across seasons. For regression models, no seasonal training is
437 required if the predictors are standardised (e.g., MLR-AAN, MLR-AAI compared to MLR-
438 RAN). Biases in the ESD methods are caused by the monthly resolution of the data. Again,
439 weather generators perform well by construction.

440 **Interannual Variability and Long-Term Trends** Interannual variability of summer
441 daily maximum temperature, measured by the variance of summer mean values, is lowest in
442 the Mediterranean and Scotland, and consistently higher in Central and Eastern Europe and
443 Scandinavia (Figure 6). ERA-Interim slightly underestimates interannual variability, again
444 likely linked to the area averaging. The performance varies widely across stations. The RCM
445 adds moderate value (high in the Mediterranean), but also spread. Simple additive MOS
446 (RaiRat-M6) leaves interannual variability unchanged. Variances of the daily distribution are
447 underestimated by ERA-Interim (see Gutiérrez and coauthors (2017)). The resulting correc-
448 tion by quantile mapping inflates interannual variability, in particular for the Mediterranean,
449 where it is overestimated by around 50%. MOS-REG underestimates interannual variability,
450 in particular when driven with ERA-Interim, because it uses predictors averaged over several
451 grid-boxes.

452 All analog methods underestimate interannual variability, consistent with the results for
453 short-term persistence. The ANALOG-ANOM method searches for continental-scale analogs
454 within a one-month window around the calendar day of interest - this likely restricts the
455 number of analogs and in turn also the represented variability. Interestingly, most regression
456 methods dramatically underestimate interannual variability. The worst performing meth-
457 ods are those without a seasonal cycle and non-standardised predictors (MLR-RAN), those
458 without temperature predictors (ESD-EOFSLP, ESD-SLP, WT-WG) and those with white
459 noise randomisation (MLR-ASW, MLR-AAW, WT-WG, SWG). Note also that both the ESD

460 methods and the SWG method are defined on continental-scale predictors, which may not be
461 suitable to capture local variations. Inflated regression by construction slightly increases the
462 variance at interannual scales. WGs do not model long-term variations and thus underestimate
463 interannual variability.

464 In addition to considering the variance at the interannual scale, we also investigate the
465 correlation between the downscaled time series and observations at the interannual scale.
466 Prior to calculating correlations, the time series are linearly detrended. This analysis provides
467 additional insight into the predictors required to explain longer-term variations. These cor-
468 relations can only be calculated when simulated and observed time series are in synchrony.
469 The RCM develops its own internal variability and thus reduces synchronicity. Therefore we
470 have not shown results for the RCM and RCM-driven MOS. Equivalently, the unconditional
471 weather generators are not in synchrony with observations and hence not shown. Correlations
472 for ERA-Interim and essentially all deterministic MOS methods are high. It is not clear to
473 the authors why CDFt and EQMWIC658 are so little synchronised - they deterministically
474 transform the ERA-Interim predictors and should thus only marginally affect the temporal
475 sequence.

476 Also PP methods perform well in general. Exceptions are the ANALOG-ANOM method,
477 the ESD methods, the WT-WG and the SWG method. Recall that ANALOG-ANOM takes
478 analogs from a 30 day window around the calendar day of interest - the identified analogs might
479 therefore have a rather strong mismatch at the local scale and thus destroy synchronicity.
480 Also, analogs of this method are defined over the whole European domain, which might result
481 in additional discrepancies at the local scales. The ESD methods, which use either 2m-
482 temperature or sea level pressure as predictor, perform worse compared to other regression
483 models; again, also the ESD method uses predictors defined over the whole of Europe. The
484 WT-WG and SWG methods perform rather bad, likely because they are based on white noise
485 randomisation. The WT-WG additionally only uses sea level pressure as predictand, the SWG
486 predictors are defined at the continental scale.

487 To characterise decadal scale variations, we considered correlations between simulated and
488 observed time series at the 7-year scale. The seasonal aggregated time series are filtered with
489 a 7-year Hamming filter. Correlations are calculated on the filtered time series without any
490 further detrending. The choice of 7 years is a compromise between the desired information
491 about long time scales, and the limited length of the time series. The effective number of data
492 points is thus low for each series (of the order of 5 per series), but still a coherent picture
493 emerges when investigating larger regions.

494 Figure ?? presents the results for summer (top panel) and winter (bottom panel) daily
495 maximum temperature. The results are overall similar to those for interannual variability.
496 Correlations are in general slightly lower during summer, in particular for ESD-SLP and
497 WT-WG (driven by sea level pressure only) for which correlations are consistently negative.
498 Correlations are lower on the Iberian Peninsula, for winter for the whole Mediterranean.

499 Finally, we investigate the representation of long-term temperature trends by the different
500 methods. Figure 8 displays the results for winter daily maximum temperatures in selected
501 regions. Of course, no results for weather generators are shown, as these do not include any
502 predictors or change factors to represent long-term changes. Note that in this experiment it is
503 not relevant whether the trends are statistically significant, because long-term variations are
504 imprinted by the ERA-Interim predictors - the right predictor choice should therefore capture
505 large-scale forced trends. It is, however, relevant whether the simulated trends are statistically

506 distinguishable from the observed trends. Thus, we calculated 95% confidence intervals of the
507 trend estimates, marked as grey shading in the panels. As trends differ very much across
508 Europe, we calculated average trends across the PRUDENCE regions. The variations of
509 trends within a region is indicated by whiskers; these denote 1.96 times the variance of all
510 trend estimates across the region.

511 Observed winter trends are highest in Scandinavia and lowest in the Mediterranean, which
512 is consistent with polar amplification. ERA-Interim performs mostly fine, but overestimates
513 trends in Central Europe, the Alps and the Mediterranean (but note that the underlying
514 ECA-D data are not homogenised, so a definite answer as to which trends are more realis-
515 tic is impossible). The RCM underestimates trends in particular in Scandinavia, but also
516 in the Alps and the Mediterranean. These trends are inherited by additive bias correction
517 (RaiRat-M6), but notably modified by many quantile mapping methods due to inflation of
518 daily variances. Note that also the ISI-MIP method, which is designed to preserve mean
519 trends, modifies trends in some regions. These trend variations are substantial, but within the
520 range of uncertainty of the observed trend estimates. The performance of PP methods again
521 depends mainly on the predictor choice. Methods using only sea level pressure or temperature
522 (but not both; ESD-EOFSLP, ESD-SLP, ESD-T2, WT-WG) tend to perform badly, although
523 filtering of stations by PCA appears to strongly increase the link with the temperature pre-
524 dictor on decadal scales (ESD-EOFT2). The ANALOG-ANOM, again, uses rather narrowly
525 defined analogs (continental scale, within one month), the SWG method combines a white-
526 noise stochastic approach with continental-scale predictors. The best performing methods
527 (ANALOG-MP, ANALOG-SP, MO-GP, MLR, MLR-WT) all include circulation predictors
528 and 2m temperature. Note, however, that 2m temperature is likely not well simulated by
529 GCMs (see the discussion in Section 4).

530 Summer trends of daily maximum temperatures (see supplementary information) are high-
531 est in Eastern Europe and the Alps. ERA-Interim in general captures these trends, but un-
532 derestimates them in the Alps and overestimates them in the Mediterranean. The RCM un-
533 derestimates summer trends everywhere, in particular in the Alps where the simulated trend
534 is not consistent with the observations. The performance of the statistical post-processing
535 methods is similar to that for winter.

536 3.2 Precipitation

537 **short-term variability** As a measure of persistence in precipitation, we consider wet-wet
538 and dry-wet transition probabilities (Figure 9). Short-term persistence in precipitation amounts
539 has not been investigated. Winter Wet-wet transition probabilities (top left panel) are low in
540 southern Europe and high along the Atlantic coasts as well as in high mountains. Winter dry-
541 wet transition probabilities (top right panel) are generally lower than wet-wet probabilities,
542 with low values in southern Europe.

543 Because it represents area average precipitation, ERA-Interim overestimates wet-wet prob-
544 abilities, in particular when further averaged. Here the RCM adds substantial value. MOS
545 methods perform consistently well. Interestingly, the simple rescaling by the method RaiRat-
546 M6 appears to perform en par with explicit wet day corrections by quantile mapping (note
547 that the BC method only treats zero precipitation as dry). MOS-AN defines analogs based on
548 simulated large-scale precipitation fields - these may not discriminate well between local dry
549 and wet days. MOS-GLM and VGLMGAMMA are both stochastic methods with white noise

550 randomisation and consequently simulate too weak wet persistence. The 4-grid-box-averaging
551 of the MOS-GLM input appears to considerably improve the performance though. Yet difficul-
552 ties in regression-based MOS techniques are evident from the low performance of MOS-GLM
553 when driven with RCM data: the RCM strongly perturbs the local day-to-day correspondence
554 between observations and simulation, which is required for a successful calibration.

555 The analog methods perform well for wet-wet transitions, most deterministic regression
556 models fail. In fact, simple linear regression models (MLR-RAN/RSN/ASW/ASI) are by
557 construction not capable of simulating daily precipitation variability - still the corresponding
558 results are included for illustration and comparison. Only the deterministic generalised linear
559 model (GLM) performs reasonably well. Most stochastic methods with white noise randomi-
560 sation (GLM-WT, WT-WG, SWG) slightly underestimate wet-day-persistence, in particular
561 WT-WG, which uses only sea level pressure, but no humidity predictors. The stochastic GLM
562 with predictors of the circulation as well as temperature and specific humidity at cloud base is
563 the best performing PP method. Interestingly, the structurally similar GLM-P (at least for the
564 occurrence process) method with similar predictors performs substantially worse. One reason
565 might be that the former defines predictors at the synoptic scale, the latter at the grid-box
566 scale. For wet-day occurrence, vertical velocities are important which can be determined from
567 horizontal convergence or divergence. Grid box pressure or velocities, however, do not carry
568 such information. Still, further analyses comparing different predictor choices are required to
569 fully understand the performance of specific predictors.

570 Dry-wet transition probabilities are well represented by ERA-Interim. The RCM has a
571 slightly positive bias. Surprisingly, however, MOS appears to reduce dry-wet transitions (by
572 wet day adjustments). Thereby it induces a negative bias for ERA-Interim, but removes the
573 positive RCM bias. Only for the UK, the positive RCM bias is even increased by many meth-
574 ods. Stochastic MOS (MOS-GLM, VGLMGAMMA) simulate too many dry-wet transitions,
575 but the averaging of simulated precipitation across grid-boxes seems to substantially improve
576 the problem (MOS-GLM-E vs. VGLMGAMMA-E). The performance of the different PP
577 methods depends strongly on both their structure and the chosen predictors. The authors
578 do not fully understand the differences in performance of different implementations. The two
579 best performing methods are ANALOG-ANOM and GLM. Both methods include circulation
580 based predictors (which should indirectly give information about lifting) and, at least indi-
581 rectly, measures of relative humidity (dew point temperature depression; specific humidity
582 in combination with temperature). Other methods, however, include similar predictors, but
583 perform worse. Recall, however, that we only know the candidate predictors used for cali-
584 bration, not the finally selected predictors at the given stations. The SS-WG and GOMEZ
585 weather generators slightly overestimate dry-wet transitions, even though this aspect is ex-
586 plicitly calibrated. Recall that the MARFI weather generator uses a wet-day threshold of 0.5
587 mm, resulting in a strong overestimation of dry-wet transitions when evaluated against a 1
588 mm threshold.

589 **Spells** The behaviour of mean spell lengths - as well as the corresponding method perfor-
590 mance - is closely tied to that of transition probabilities (Figure 10). Mean winter wet-spell
591 lengths (top left) are high along the Atlantic west coasts and mountain ranges, and
592 short in Eastern Europe and the Mediterranean. Summer dry spells (top right) are short in
593 Central and Northern Europe, and long in the Mediterranean.

594 ERA-Interim underestimates winter wet spells because of spatial averaging (upper panel).

595 At first sight, the RCM adds no value. Yet the RCM reduces the ERA-Interim bias of too
596 many wet-days Gutiérrez and coauthors (2017) as well as the bias in too high a wet-wet
597 transition probability (see above). As a result, the RCM implicitly adds value in the sub-
598 sequent bias correction, in particular over the Iberian Peninsula. Quantile mapping without
599 seasonal training (GQM, GPQM, EQM) overestimates winter wet spell lengths. Interestingly,
600 conditioning on weather types (EQM-WT) essentially has the same effect as an explicit sea-
601 sonal training (EQMs), indicating that biases are circulation dependent and translate into
602 seasonally-dependent biases, because the frequency of weather types changes throughout the
603 year. The MOS-AN, MOS-GLM and VGLMGAMMA perform very similar as with regard to
604 short-term persistence. In particular the averaging of predictors across 4 grid boxes in the
605 stochastic methods (MOS-GLM-E vs. VGLMGAMMA-E) seems to be crucial to increase skill.
606 The performance of the PP methods scatters widely, as already for short term persistence.
607 Only the ANALOG-ANOM and GLM perform well. The SS-WG and GOMEZ Weather gen-
608 erators slightly underestimate wet spell lengths. Again, the MARFI weather generator sticks
609 out because of the different wet day threshold.

610 The performance for summer dry spells is overall similar to that for winter wet spells.
611 ERA-Interim spells are again too short, but here the RCM adds substantial value, likely due
612 to a reduction of the area-average-related drizzle effect. MOS appears to increase the length of
613 dry spells as a consequence of the wet day correction. For ERA-Interim this leads to unbiased
614 results, whereas the RCM performance is deteriorated towards too long dry spells. This
615 problem occurs in particular for quantile mapping methods, which are not seasonally trained
616 (GQM, GPQM, EQMs, EQM-WT). Analog methods perform slightly better for dry- than for
617 wet spells, the GLM performs worse than for wet spells, but still reasonably well. Weather
618 generators perform slightly better for dry- than for wet spells. Owing to the different wet-day
619 threshold, the MARFI weather generator is slightly more biased and has a much higher spread
620 across stations. In general, the length of dry spells is overestimated in the Mediterranean and
621 France.

622 **Seasonality** Seasonality of precipitation is measured by the relative amplitude (defined
623 as the difference between precipitation in the maximum and minimum of the seasonal cycle,
624 relative to the annual mean) and phase (defined as the position of the maximum of the seasonal
625 cycle). Although the calculation is identical to that of the seasonal cycle of temperature,
626 some details will be relevant in particular for precipitation. In fact, the seasonal cycle of
627 precipitation has two peaks in many regions, sometimes even shoulders or peaks that may be
628 artefacts of sampling variability. Following Favre et al. (2016), we therefore filter the seasonal
629 cycle by four harmonics - this model is flexible enough to capture smooth - likely physical -
630 variations, but at the same time filters out residual noise (see Figure 2). The amplitude of
631 the seasonal cycle is simply defined as the difference between maximum and minimum. For
632 the phase definition, further steps have been carried out. They are a compromise between
633 being simple and transparent, but at the same time capturing the complex seasonal behaviour.
634 First, secondary peaks with an amplitude (defined as the difference between the closest local
635 minimum and the peak itself) of less than 10% of the total amplitude have been removed,
636 as well as neighboring peaks with a minimum in between that is less than 10% of the total
637 amplitude lower than the mean height of the two peaks. The two peaks are then replaced
638 by a single peak by averaging their height as well as phase. The first step removes all minor
639 peaks, the second step removes dips in an overall broad maximum, which are both likely an

640 artefact of sampling variability. Visual inspection of observed seasonality for all 86 stations
641 corroborates that this definition conforms with expert judgment. We then record the phase
642 of the remaining highest and second highest peak for observations and all simulations. The
643 observed phase is then defined as that of the highest peak. The simulated phase is defined as
644 the phase of that of the two highest peaks, which is closest to the observed. The latter definition
645 avoids that, if highest and second highest peak have similar height and are swapped in the
646 simulation, an artificially large phase bias is calculated. Apart from this phase definition we
647 considered other measures for characterising the timing of the seasonal cycle, but rejected all
648 other possibilities. We considered, e.g., correlations between simulated and observed seasonal
649 cycle, but this measure is difficult to interpret in terms of an actual mismatch in timing.
650 Additionally, we also considered to calculate phases of secondary peaks, but concluded that a
651 plain and transparent presentation of performance across Europe would be difficult.

652 Seasonality of precipitation (Figure 11) has a strong north-south gradient, ranging from
653 less than 50% of annual mean precipitation in central-west Europe to more than 200% in
654 southern Spain and southern Greece. The annual cycle peaks in winter along the Atlantic and
655 the Mediterranean, and in summer in Central and eastern Europe and eastern Scandinavia.
656 Reanalysis and RCM underestimate the amplitude of the annual cycle, although the RCM
657 adds considerable value. MOS generally performs well, although methods without seasonal
658 training (GQM, GPQM, EQM, EQM-WT) overestimate the relative amplitude by about 20%.
659 Note, however, that conditioning the correction on weather types (EQM-WT) substantially
660 reduces this bias. PP performance again depends on the method-type, the treatment of
661 seasonality, and the choice of predictors. The analog methods perform reasonably well, linear
662 regression models all underrepresent the relative amplitude (MLR-RAN/RSN/ASW/ASI).
663 The good performance of the GLM method indicates that a sensible model structure and
664 predictor choice (circulation and humidity) may allow to capture the seasonal cycle without
665 an explicit model. The phase of the seasonal cycle is well captured by most methods. The
666 bad performance of WT-WG indicates that sea level pressure alone does not determine the
667 seasonal cycle.

668 **Interannual Variability and long-term trends** Interannual variability of precipitation
669 varies unsystematically in space (Figure 12). Values, however, tend to be higher at higher ele-
670 vations. As for temperature, reanalysis data underrepresent interannual variability, especially
671 at low resolution. But in contrast to temperature, the RCM succeeds in reducing the overall
672 bias, in particular over the Mediterranean. Deterministic MOS methods suffer strongly from
673 variance inflation, which in cases doubles the interannual variance. Regression based MOS by
674 contrast tends to underestimate interannual variability, consistent with the driving model.
675 The performance of PP methods, again, varies considerably. Note, however, that all well
676 performing methods include not only circulation-based predictors, but also measures of hu-
677 midity (ANALOG-ANOM, ANALOG, ANALOG-SP, GLM-det, GLM, GLM-WT). Weather
678 generators, as expected, underestimate interannual variability - even more so for the MARFI
679 weather generator because of the different wet-day threshold.

680 Interannual correlations are, as expected, lower for precipitation than for temperature:
681 only about 50% of the local variability ($\sim 0.7^2$) seems to be explained by the area average,
682 the rest is due to local variability. Deterministic MOS methods do not modify this correlation
683 (again, we cannot explain the performance of EQM-WIC658). For the stochastic MOS meth-
684 ods, the value of averaging simulated precipitation across neighboring grid boxes is evident

685 (compare MOS-GLM-E and VGLMGAMMA-E). All PP methods explain substantially less
686 of the interannual variability than the grid-box ERA-Interim. The worst performing methods
687 are ANALOG-ANOM (analogs searched within 30 day window only, continental scale pre-
688 dictors and analogs), MLR-ASW (Gaussian white noise radomisation), WT-WG (stochastic,
689 only sea level pressure as predictors) and SWG (stochastic, continental scale predictors). Note
690 the substantial difference between the - structurally similar - GLM and SWG models. GLM
691 defines predictors on a national scale, SWG on a continental scale.

692 Seven year correlations between simulations and observations are similar to interannual
693 correlations; they are much higher though in winter than in summer (see supplementary
694 information).

695 Finally, we investigate the performance in representing relative trends in seasonal mean
696 precipitation. Figure 13 presents the results for summer and selected regions. All observed
697 trends are essentially zero and insignificant, with moderately positive values in Central Europe.
698 We nevertheless show the results to demonstrate the behaviour of the different methods. ERA-
699 Interim captures the observed trends in some regions, but simulates a zero trend for Central
700 Europe, and a negative trend for the Alps. The RCM simulates positive trends for the British
701 Isles, Central Europe, Scandinavia and the Alps, although all these are within the range of
702 sampling uncertainty. The MOS methods tend to inflate the wrong RCM trends, as well as the
703 wrong negative ERA-Interim trends in the Alps. Many PP methods capture observed trends
704 quite well, although the performance changes substantially - and not for obvious reasons - from
705 region to region. Identifying necessary predictors appears to be much less straight forward than
706 in case of temperature trends.

707 4 Discussion and Conclusions

708 We have systematically evaluated how different types of downscaling and bias correction ap-
709 proaches represent temporal aspects. These aspects comprise systematic seasonal variations
710 and residual temporal dependence such as short-term persistence, spell length distributions
711 and interannual to long-term variability variability. Additionally, we considered long-term
712 trends, which are a superposition of long-term internal climate variability and forced trends.
713 Our results complement, corroborate and extend earlier findings, in particular by Frost et al.
714 (2011), Hu et al. (2013), Benestad and Haugen (2007) and Huth et al. (2015).

715 Overall, the behaviour of the different approaches turned out to be as expected from their
716 structure and implementation. For the interpretation of the results, it has to be acknowledged
717 whether a particular aspect of a model is explicitly calibrated - a good performance is then
718 more or less trivial - or emerges from the model, e.g., by well chosen meteorological predictors.

719 A summary of the results (apart from correlations and long-term trends) can be found in
720 Figure 14. The raw ERA-Interim data are typically biased compared to observed station data,
721 stronger so for the spatially aggregated 2° version. Note, however, that these discrepancies are
722 not necessarily bias in the sense of model errors, but simply reflect the scale-gap between area
723 averages and point values (Volosciuk et al., 2015). The chosen RCM adds value to reanalysis
724 data for most considered aspects, for all seasons and for both temperature and precipitation.
725 Note, however, that we included just one RCM in our validation study. One should be careful
726 in generalising these results because RCMs may differ considerably in their ability to reproduce
727 temporal characteristics (Kotlarski et al., 2014; Huth et al., 2015).

728 The MOS methods considered in this intercomparison do not explicitly change the resid-
729 ual temporal dependence (and it is questionable whether they should explicitly do so, as such
730 changes would destroy the temporal consistency with the driving model). However, quantile
731 mapping approaches modifying the marginal distribution (including wet day probabilities)
732 do indirectly improve temporal variability. For temperature, some implementations slightly
733 improve short-term persistence, but in particular for precipitation, the representation of tran-
734 sition probabilities as well as wet and dry spells is substantially improved. Interestingly,
735 dry-wet transitions and dry-spell lengths are much better for the bias-corrected RCM than
736 for bias-corrected reanalyses, even though the added value of the RCM for these indices was
737 marginal only. Interannual and long-term variability is typically inflated by MOS. Moder-
738 ately for temperature, but substantially for precipitation. These findings corroborate earlier
739 results of adverse inflation effects by quantile mapping (Maraun, 2013). long-term trends are
740 inherited from the driving model, but may be substantially deteriorated by further variance
741 inflation. The annual cycle is improved by almost all MOS methods - but recall that most
742 methods are seasonally trained. Conditioning on weather types (EQM-WT) seems to a suc-
743 cessful - and physically more defensible - variant to better represent the annual cycle. In
744 any case, our results clearly show that - for many but not all temporal aspects - dynamical
745 downscaling prior to the bias correction substantially improves the results compared to a di-
746 rect bias correction from the global model⁵. The reason of course is that the bias correction
747 does not improve the representation of meso-scale processes. Thus, depending on the context,
748 dynamical downscaling may be advisable or even essential.

749 The performance of the participating PP methods varies strongly from aspect to aspect
750 and method to method. Analogue methods show difficulties representing temperature vari-
751 ability, but perform quite well for precipitation variability. Two reasons may contribute to
752 the low performance for temperature: first, predictors describing circulation and humidity
753 have much stronger spatial-temporal variability than temperature fields and therefore domi-
754 nate the definition of the analogs. Second, predictors and analogs are often defined on large
755 scales. Locally, differences between actual weather and analogs may be substantial. Thus,
756 even if analogs may describe a smooth temperature evolution at large scales, the resulting
757 local sequence might be too noisy.

758 Deterministic linear regression models perform fairly well for temperature, but overesti-
759 mate short-term persistence and spell lengths. White noise randomisation deteriorates the
760 representation of these aspects. Linear regression models, in any variant, are far too sim-
761 plistic for precipitation downscaling. They strongly overestimate wet-wet transitions and the
762 length of wet spells, while stochastic methods underestimate these aspects. Biases for dry-
763 wet transitions and dry-spell lengths tends to be opposite to those for wet-wet transitions
764 and wet-spell lengths, but they are substantial for almost all PP methods. Only a stochastic
765 generalised linear model with suitable predictors has shown to perform well (GLM). A struc-
766 tually similar model (SWG) - with similar predictor variables, but defined on the continental
767 scale - performs notably bad. The representation of the annual cycle depends strongly on
768 the individual method; whether or not a method is seasonally trained plays a minor role -
769 the choice of reasonable predictors seems to be a key factor. For temperature, temperature
770 related predictors are required; for precipitation, circulation and humidity based predictors.
771 There is evidence that biases in interannual variability of temperature mainly depend on the

⁵Note in this context, that the ERA-Interim is an “ideal” GCM in the sense that it is forced to closely follow the observed large-scale weather.

772 method type (again, analog methods and white noise randomisation underestimate internal
773 variability), on the predictor variables (all well performing methods combine circulation and
774 temperature predictors) and the domain size (all methods using continental-size predictor
775 domains perform badly). For precipitation, the inclusion of predictors that represent both
776 circulation and humidity appears to be crucial. long-term trends in temperature are cap-
777 tured by models with surface temperature predictors (see the critical discussion below), for
778 precipitation no conclusions can be drawn based on the available ensemble, and the rather
779 low signal-to-noise ratio. Overall, white-noise randomisation with continental-scale predictors
780 turned out to perform weakly. Apparently, the variance explained by predictors at such large
781 scales is rather low, such that the residual white noise is too strong to retain the overall
782 temporal dependence.

783 Unconditional weather generators tend to perform well for the aspects they have been
784 calibrated for: they only slightly underestimate short-term temperature persistence and wet-
785 wet transitions, but slightly overestimate dry-wet transitions. Nevertheless also many non-
786 calibrated aspects are fairly well represented. Temperature spell lengths are slightly underes-
787 timated, in particular for winter cold spells and summer warm spells. Wet spell lengths are
788 well represented, dry spell lengths underestimated. Only interannual variability is substan-
789 tially underrepresented. These effects are well known issues (Wilks and Wilby, 1999) and are
790 relevant also for decadal variability. Seasonality is, by construction, well simulated.

791 Overall, the performance is similar in different seasons - but recall that in particular most
792 MOS methods and all weather generators are calibrated to do so. These explicit seasonal
793 models, however, may be questioned for being used in a future climate: seasonally varying
794 biases indicate that seasonal biases may also change differently on long time scales.

795 Our findings highlight a series of open research questions, and the need for a range of
796 improvements. MOS methods perform overall very well. Some key issues, however, remain
797 to be addressed: the inflation (or potentially deflation) of interannual and long-term vari-
798 ability and trends is of course directly tied to the simplicity of quantile mapping compared
799 to MOS methods in weather forecasting and the PP methods presented here: whereas the
800 latter express physical relationships between large and local scales at least rudimentarily as
801 regression models and thereby can distinguish between forced and local internal variability,
802 quantile mapping adjusts only long-term distributions of daily values without any physical
803 basis. This calibration is especially problematic when a scale gap between predictand and
804 predictor is to be bridged (Maraun, 2013). The reason for the calibration, of course, is that
805 regression models cannot easily be calibrated in a free running climate model, which is not in
806 synchrony with observations Maraun et al. (2010). More research is needed to understand the
807 link between biases in short-term variability and long-term variability. Some methods have
808 been developed to separate variability on different scales, and to adjust them independently,
809 other methods have been developed to preserve climate model trends to various degrees (Li
810 et al., 2010; Haerter et al., 2011; Hempel et al., 2013; Pierce et al., 2015). The physical as-
811 sumptions underlying these different methods need to be better understood. In any case, our
812 results show that any bias correction relies on climate models that simulate realistic trends.
813 In case of downscaling to a finer resolution, it might be useful to separate the bias correction
814 from the downscaling, i.e., apply a correction against gridded observational data, and then
815 implement a stochastic downscaling model against point data (Volosciuk et al., 2017). Re-
816 gression based MOS methods have been presented as further alternatives (MOS-REG/GLM,
817 VGLMGAMMA), but these cannot be calibrated to standard climate model simulations. The

818 results show that even typical RCM hindcast simulations (where the RCM is driven with a
819 reanalysis, MOS-REG-R and MOS-GLM-R) are not sufficiently synchronous to ensure a suc-
820 cessful calibration. A way out might be to condition bias correction on weather types, such
821 as demonstrated by EQM-WT.

822 Various research strands are possible and necessary to better understand and to improve
823 PP methods. For analog methods, in particular in case of temperature, a way forward could
824 be based on defining the analogs not on a single day, but rather on a sequence of days (e.g.
825 Beersma and Buishand, 2003). Such approaches, however, require long time series. Note, how-
826 ever, that analog methods cannot represent substantial climatic changes, where no analogs
827 might be available to sample from Gutiérrez et al. (2013). An obvious improvement of regres-
828 sion models is a better representation of residual variability - for temperature the in linear
829 models for temperature, and generalised linear models for precipitation. Here, conditional
830 weather generators are promising that extend the white noise randomisation (both for tem-
831 perature and precipitation) by a Markov component. For instance, one may include not only
832 meteorological predictors, but also simulated predictand values from previous days as predic-
833 tors (Chandler and Wheeler, 2002; Yang et al., 2005).

834 The crucial questions regarding the PP approach are, however, not an improvement in
835 model structure, but a better understanding of predictor choice. Unfortunately, the available
836 model ensemble did not allow for a stringent identification of suitable predictors. Nevertheless,
837 the results highlight a couple of issues. Note that these are questions of physics more than of
838 statistics. First, what is a suitable domain size? The GLM-P and GLM methods include a
839 structurally similar rainfall occurrence process and a - at first sight - similar set of predictors.
840 But the GLM method performs far better than GLM-P in simulating all occurrence-related
841 aspects. A major difference between the two implementations is that GLM uses synoptic
842 scale predictors, whereas GLM-P relies on grid-box predictors. Precipitation occurrence is
843 controlled by relative humidity and vertical velocity. The latter is typically represented by
844 predictors of the horizontal circulation. The underlying reasoning is that horizontal divergence
845 and convergence determines vertical descent and ascent. Convergence and divergence, in
846 turn, may be implicit in large-scale pressure fields, but they are not represented by grid-box
847 pressure values. Thus, the choice of predictor variables depends on the domain size. Many
848 methods with limited performance, in particular for temperature, where based on continental-
849 scale predictors. Thus, there is evidence that such predictor domains are simply too large to
850 successfully represent local variability. Here one has to trade-off between downscaling across
851 large areas and precision at local scales. In fact, we see the main strength of PP methods not in
852 competing with RCMs across whole continents, but rather in providing tailored region-specific
853 projections.

854 Second, which predictors are required for representing long-term trends? We demonstrated
855 that model performance for the same set of predictors differed substantially for short-term per-
856 sistence and long-term changes. The reason of course is that downscaling methods are cali-
857 brated to day-to-day-variability, but are intended to work on long-term variability (Huth et al.,
858 2015). For temperature, a combination of temperature and circulation predictors appeared
859 to fail well explain long-term trends. Precipitation, however, is a more complex nonlinear
860 process, and no method convincingly captured trends in all considered regions. A further
861 complicating issue is the low signal to noise ratio: all trends, and all misrepresentations, are
862 still within the sampling uncertainty.

863 Weather generators do have an explicit model of the short-term temporal dependence,

864 but those variants participating in this intercomparison did not include any meteorological
865 predictors. As a result, these methods underestimated long-term variability - it was not
866 explicitly modelled. Also here improvements are possible, e.g., by conditioning the weather
867 generator on monthly aggregates (being generated by the separate monthly WG or taken from
868 the driving data - e.g. GCM, RCM or reanalysis) to improve the representation of interannual
869 variability (Dubrovský et al., 2004).

870 This study was based on a perfect predictor setting to isolate downscaling skill. Therefore,
871 we did not investigate the performance with imperfect predictors or boundary conditions from
872 free running GCMs. Downscaling methods - apart from unconditional weather generators -
873 to a large extent inherit the errors in representing temporal variability of the driving models
874 (Hall, 2014). The downscaling performance may, therefore, drop considerably, when driven by
875 imperfect forcing from a GCM. For MOS, the issue is rather subtle: marginal biases in present
876 climate are by construction removed, hence it is difficult to identify fundamental GCM errors
877 such as the misrepresentation of the large-scale circulation and its temporal structure. Thus,
878 also non-calibrated aspects, in particular the temporal aspects, should thus be evaluated.

879 For PP one typically assumes that large-scale predictors from the free atmosphere fulfill
880 the PP assumption. This assumption should be tested for GCMs. Again, evaluating temporal
881 aspects might be more informative than evaluating marginal aspects - often, predictors are
882 based on anomalies, such that mean biases are implicitly removed. But even more, many PP
883 predictors are not defined at large scales, and not chosen from the free atmosphere. For in-
884 stance, those methods that best represented temperature trends all relied on 2m-temperature.
885 In the reanalysis, which has been used as predictors, temperature observations have been
886 assimilated into the model, such that grid-box variability and long-term are likely correctly
887 represented in data rich regions. Local surface feedbacks that modulate temperature vari-
888 ability are thus implicitly accounted for. But a free running GCM will likely not correctly
889 represent these feedbacks, such that GCM simulated 2m temperature will likely not fulfill the
890 PP assumption. Similar arguments apply for grid box values of, e.g., 10m winds.

891 Even though we investigated the performance to represent observed trends, we can only
892 draw limited conclusions about representing future trends. MOS relies on credibly simulated
893 grid box trends - the ERA-Interim trends are approximately correct by construction, the
894 RCM show substantial deficiencies. But also for PP methods, our findings are far from being
895 conclusive. For temperature, as discussed before, the PP assumption for relevant predictors
896 may not be fulfilled. For precipitation, simply no conclusions are possible because of the low
897 signal-to-noise ration. In any case, a method performing badly with perfect predictors will not
898 perform better with imperfect predictors. Passing this evaluation is therefore a necessary, but
899 not a sufficient requirement for a method to be applicable under climate change conditions.

900 This discussion shows that further studies are required to establish the skill of down-
901 scaling under simulated future conditions. The VALUE community is planning additional
902 experiments Maraun et al. (2015): GCM predictor experiments to asses the performance un-
903 der imperfect predictors, and pseudo reality experiments to establish statistical downscaling
904 skill in simulated future climates. Additionally, we have identified a range of open questions
905 that can be addressed within our perfect predictor experiment, in particular related to the
906 predictor choice of PP methods. The metadata and complete results for individual methods
907 are available from the VALUE portal www.value-cost.eu/validationportal. They can be
908 downloaded and further analysed. Additionally, we encourage dedicated sensitivity studies
909 based on the ensemble at hand.

910 Appendix

911 Similarly to the portrait diagram in Sillmann et al. (2013), Figure 14 summarises the perfor-
912 mance of the different methods for different indices in one (color-coded) value. To make these
913 comparable across methods and indices, a reference scale has to be defined. This scale cannot
914 simply be measured in terms of the best and worst performing methods for an index, as such
915 a scale would only measure relative performance, not absolute performance. For instance, one
916 would not be able to distinguish an index that is well represented from one that is poorly
917 represented by all methods. Sillmann et al. (2013) define the variability of an index in space
918 as reference scale. But this scale cannot be applied to a single series, and it cannot distinguish
919 between indices that are well modelled by all methods across space (e.g., the seasonal cycle)
920 and indices that are badly modelled (e.g., interannual variability). Thus, we attempt to define
921 natural scales for different types of indices:

- 922 • For biases in mean temperature, we define twice the standard deviation of daily vari-
923 ability as scale. For Gaussian distributed variables, this range spans roughly 95% of the
924 probability mass.
- 925 • For biases of temperature indices, which may be expressed as anomalies (such as the 20
926 year return value or the amplitude of the seasonal cycle), we chose the actual modulus
927 of the anomaly (i.e., the difference of the return value and mean temperature, or the
928 amplitude itself) as reference scale.
- 929 • For relative biases, which assume only positive values (such as for temperature variance,
930 precipitation intensity or mean spell length), a natural scale is the observed value itself.
- 931 • For the phase of the seasonal scale we (somewhat arbitrarily) define one month as a
932 reference scale.

933 References

- 934 A. Alexandru, R. de Elia, and R. Laprise. Internal variability in regional climate downscaling
935 at the seasonal scale. *Mon. Wea. Rev.*, 135(9):3221–3238, 2007. doi: 10.1175/MWR3456.1.
- 936 J. Bartholy, R. Pongrácz, and A. Kis. Projected changes of extreme precipitation using multi-
937 model approach. *Q.J. Hung. Meteorol. Serv.*, 119:129–142, 2015.
- 938 J. Bedia, M. Iturbide, S. Herrera, R. Manzanas, and J. Gutiérrez. down-
939 scaler: Climate data manipulation, bias correction and statistical downscaling.
940 <http://github.com/SantanderMetGroup/downscaleR/wiki>, 2016.
- 941 J.J. Beersma and T.A. Buishand. Multi-site simulation of daily precipitation and temperature
942 conditional on the atmospheric circulation. *Clim. Res.*, 25:121–133, 2003.
- 943 R. Benestad, A. Mezghani, and K. Parding. esd: Climate analysis and
944 empirical-statistical downscaling (ESD) package for monthly and daily data.
945 <http://rcg.gvc.gu.se/edu/esd.pdf>, 2015a.

- 946 R.E. Benestad and J.E. Haugen. On complex extremes: flood hazards and combined high
947 spring-time precipitation and temperature in Norway. *Climatic Change*, 85(3-4):381–406,
948 2007.
- 949 R.E. Benestad, D. Chen, A. Mezghani, L. Fan, and K. Parding. On using principal components
950 to represent stations in empirical-statistical downscaling. *Tellus A*, 67, 2015b.
- 951 M.S. Bukovsky. Temperature trends in the NARCCAP regional climate models. *J. Climate*,
952 24:3985–3991, 2012.
- 953 G. Bürger, T. Q. Murdock, A. T. Werner, S. R. Sobie, and A. J. Cannon. Downscaling extremes
954 - an intercomparison of multiple statistical methods for present climate. *J. Climate*, 25:
955 4366–4388, 2012.
- 956 P. Calanca. Climate change and drought occurrence in the Alpine region: How severe are
957 becoming the extremes? *Glob. Planet. Change*, 57(1):151–160, 2007.
- 958 A.J. Cannon. Multivariate Bias Correction of Climate Model Output: Matching Marginal
959 Distributions and Intermittent Dependence Structure. *J. Climate*, 29(19):7045–7064, 2016.
- 960 P. Ceppi, S.C. Scherrer, A.M. Fischer, and C. Appenzeller. Revisiting Swiss temperature
961 trends 1959–2008. *Int. J. Climatol.*, 32:203–213, 2012.
- 962 R. E. Chandler and H. S. Wheater. Analysis of rainfall variability using generalized linear
963 models: A case study from the west of Ireland. *Wat. Resour. Res.*, 38(10):1192, 2002.
- 964 S.P. Charles, B.C. Bates, and J.P. Hughes. A spatiotemporal model for downscaling precipi-
965 tation occurrence and amounts. *J. Geophys. Res.*, 104(D24):31,657–31,669, 1999.
- 966 J. H. Christensen and O. B. Christensen. A summary of the PRUDENCE model projections
967 of changes in European climate by the end of this century. *Clim. Change*, 81:7–30, 2007.
- 968 D.P. Dee, S.M. Uppala, A.J. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae, M.A.
969 Balmaseda, G. Balsamo, P. Bauer, P. Bechtold, A.C.M. Beljars, L. van den Berg, J. Bidlot,
970 N. Bormann, C. Delsol, R. Dragani, M. Fuentes, A.J. Geer, L. Haimberger, S.B. Healy,
971 H. Hersbach, E.V. Hólm, L. Isaksen, P. Kållberg, M. Köhler, M. Matricardi, A.P. McNally,
972 B.M. Monge-Sanz, J.-J. Morcrette, B.-K. Park, C. Peubey, P. de Rosnay, C. Tavolato, J.-N.
973 Thépaut, and F. Vitart. The ERA-Interim reanalysis: configuration and performance of
974 the data assimilation system. *Quart. J. Royal Meteorol. Soc.*, 137:553–597, 2011.
- 975 M. Dubrovský, J. Buchtele, and Z. Žalud. High-frequency and low-frequency variability in
976 stochastic daily weather generator and its effect on agricultural and hydrologic modelling.
977 *Clim. Change*, 63(1):145–179, 2004.
- 978 A. Favre, N. Philippon, B. Pohl, E.-A. Kalognomou, C. Lennard, B. Hewitson, G. Nikulin,
979 A. Dosio, H.-J. Panitz, and R. Cerezo-Mota. Spatial distribution of precipitation annual
980 cycles over South Africa in 10 CORDEX regional climate model present-day simulations.
981 *Clim. Dynam.*, 46:1799–1818, 2016.
- 982 E.M. Fischer, S.I. Seneviratne, P.L. Vidale, D. Lüthi, and C. Schär. Soil moisture-atmosphere
983 interactions during the 2003 European summer heat wave. *J. Climate*, 20:5081–5099, 2007.

- 984 H. J. Fowler, S. Blenkinsop, and C. Tebaldi. Linking climate change modelling to impacts stud-
985 ies: recent advances in downscaling techniques for hydrological modelling. *Int. J. Climatol.*,
986 27:1547–1578, 2007.
- 987 C. Frei, J. H. Christensen, M. Deque, D. Jacob, R. G. Jones, and P. L. Vidale. Daily precipita-
988 tion statistics in regional climate models: Evaluation and intercomparison for the european
989 alps. *Journal of Geophysical Research-Atmospheres*, 108(D3), 2003.
- 990 P. Froidevaux, J. Schwanbeck, R. Weingartner, C. Chevalier, and O. Martius. Flood triggering
991 in Switzerland: the role of daily to monthly preceding precipitation. *Hydrol. Earth Syst.*
992 *Sci.*, 19(9):3903–3924, 2015.
- 993 A.J. Frost, S.P. Charles, B. Timbal, F.H.S. Chiew, R. Mehrotra, K.C. Nguyen, R.E. Chandler,
994 J.L. McGregor, G. Fu, D.G.C. Kirono, et al. A comparison of multi-site daily rainfall
995 downscaling techniques under Australian conditions. *J. Hydrol.*, 408(1):1–18, 2011.
- 996 F. Giorgi, X. Bi, and J. Pal. Mean, interannual variability and trends in a regional climate
997 change experiment over Europe. I. Present-day climate (19611990). *Clim. Dynam.*, 22:
998 733–756, 2004.
- 999 C.M. Goodess, C. Anagnostopoulou, A. Bárdossy, C. Frei, C. Harpham, M.R. Haylock,
1000 Y. Hundecha, P. Maheras, J. Ribalaygua, J. Schmidli, T. Schmith, K. Tolika, R. Tomozeiu,
1001 and R.L. Wilby. An intercomparison of statistical downscaling methods for Europe and
1002 European regions assessing their performance with respect to extreme weather events and
1003 the implications for climate change applications. Technical report, Climatic Research Unit,
1004 2010.
- 1005 J.M. Gutiérrez and coauthors. An intercomparison of a large ensemble of statistical down-
1006 scaling methods for europe: Overall results from the value perfect predictor cross-validation
1007 experiment. *Int. J. Climatol.*, *subm.*, 2017.
- 1008 J.M. Gutiérrez, D. San-Martín, S. Brands, R. Manzanas, and S. Herrera. Reassessing statistical
1009 downscaling techniques for their robust application under climate change conditions. *J.*
1010 *Climate*, 26(1):171–188, 2013.
- 1011 E. Gutmann, T. Pruitt, M.P. Clark, L. Brekke, J.R. Arnold, D.A. Raff, and R.M. Rasmussen.
1012 An intercomparison of statistical downscaling methods used for water resource assessments
1013 in the united states. *Wat. Resour. Res.*, 50(9):7167–7186, 2014.
- 1014 J.O. Haerter, S. Hagemann, C. Moseley, and C. Piani. Climate model bias correction and the
1015 role of timescales. *Hydrol. Earth Syst. Sci.*, 15(3):1065–1079, 2011.
- 1016 A. Hall, X. Qu, and J.D. Neelin. Improving predictions of summer climate change in the
1017 united states. *Geophys. Res. Lett.*, 35:L01702, 2008.
- 1018 Alex Hall. Projecting regional change. *Science*, 346(6216):1461–1462, 2014.
- 1019 M. R. Haylock, G. C. Gawley, C. Harpham, R. L. Wilby, and C. M. Goodess. Downscaling
1020 heavy precipitation over the United Kingdom: A comparison of dynamical and statistical
1021 methods and their future scenarios. *Int. J. Climatol.*, 26(10):1397–1415, 2006.

- 1022 S. Hempel, K. Frieler, L. Warszawski, J. Schewe, and F. Piontek. A trend-preserving bias
1023 correction - the ISI-MIP approach. *Earth Syst. Dynam.*, 4:219–236, 2013.
- 1024 S. Herrera, M. Turcu, and J.M. Gutiérrez. A mos-regression technique for temporally-coherent
1025 bias correction of regional climate model simulations. *Clim. Dynam.*, submitted, 2017.
- 1026 E. Hertig and coauthors. Validation of extremes from the perfect-predictor experiment of the
1027 cost action value. *Int. J. Climatol.*, 2016.
- 1028 E. Hertig and J. Jacobeit. Assessments of Mediterranean precipitation changes for the 21st
1029 century using statistical downscaling techniques. *Int. J. Climatol.*, 28:1025–1045, 2008.
- 1030 Y. Hu, S. Maskey, and S. Uhlenbrook. Downscaling daily precipitation over the yellow river
1031 source region in china: a comparison of three statistical downscaling methods. *Theor. Appl.*
1032 *Climatol.*, 112(3-4):447–460, 2013.
- 1033 R. Huth. Statistical downscaling of daily temperature in central europe. *J. Climate*, 15:
1034 1731–1742, 2002.
- 1035 R. Huth, J. Kyselý, and M. Dubrovský. Time structure of observed, GCM-simulated, down-
1036 scaled, and stochastically generated daily temperature series. *J. Climate*, 14:4047–4061,
1037 2001.
- 1038 R. Huth, J. Miksovsky, P. Stepanek, M. Belda, A. Farda, Z. Chladova, and P. Pisoft. Compar-
1039 ative validation of statistical and dynamical downscaling models on a dense grid in central
1040 Europe: temperature. *Theor. Appl. Climatol.*, 120(3-4):533–553, MAY 2015. ISSN 0177-
1041 798X. doi: 10.1007/s00704-014-1190-3.
- 1042 D. Jacob, L. Bärring, O. B. Christensen, J. H. Christensen, M. de Castro, M. Déqué, F. Giorgi,
1043 S. Hagemann, M. Hirschi, R. Jones, E. Kjellström, G. Lenderink, B. Rockel, E. Sánchez,
1044 C. Schär, S. I. Seneviratne, S. Somot, A. van Ulden, and B. van den Hurk. An inter-
1045 comparison of regional climate models for Europe: model performance in present-day cli-
1046 mate. *Clim. Change*, 81:31–52, 2007.
- 1047 E.-A. Kalognomou, C. Lennard, M. Shongwe, I. Pinto, A. Favre, M. Kent, B. Hewitson,
1048 A. Dosio, G. Nikulin, H.-J. Panitz, and M. Büchner. A Diagnostic Evaluation of Precipita-
1049 tion in CORDEX Models over Southern Africa. *J. Climate*, 26(23):9477–9506, 2013. doi:
1050 10.1175/JCLI-D-12-00703.1.
- 1051 D. Keller, A.M. Fischer, C. Frei, M.A. Liniger, C. Appenzeller, and R. Knutti. Implementation
1052 and validation of a Wilks-type multi-site daily precipitation generator over a typical Alpine
1053 river catchment. *Hydrol. Earth Syst. Sci.*, 19:2163–2177, 2015.
- 1054 D.E. Keller, A.M. Fischer, M.A. Liniger, C. Appenzeller, and R. Knutti. Testing a weather
1055 generator for downscaling climate change projections over Switzerland. *Int. J. Climatol.*,
1056 2016.
- 1057 C. G. Kilsby, P. D. Jones, A. Burton, A. C. Ford, H. J. Fowler, C. Harpham, P. James,
1058 A. Smith, and R. L. Wilby. A daily weather generator for use in climate change studies.
1059 *Env. Mod. Soft.*, 22:1705–1719, 2007.

- 1060 A.M.G. Klein Tank, J.B. Wijngaard, G.P. Können, R. Böhm, G. Demarée, A. Gocheva,
1061 M. Mileta, S. Pashiardis, L. Hejkrlik, C. Kern-Hansen, R. Heino, P. Bessemoulin, G. Müller-
1062 Westermeier, M. Tzanakou, S. Szalai, T. Pálsdóttir, D. Fitzgerald, S. Rubin, M. Capaldo,
1063 M. Maugeri, A. Leitass, A. Bukantis, R. Aberfeld, A.F.V. van Engelen, E. Forland, M. Mi-
1064 etus, F. Coelho, C. Mares, V. Razuvaev, E. Nieplova, T. Cegnar, J.A. López, B. Dahlström,
1065 A. Moberg, W. Kirchhofer, A. Ceylan, O. Pachaliuk, L.V. Alexander, and P. Petrovic. Daily
1066 dataset of 20th-century surface air temperature and precipitation series for the European
1067 climate assessment. *Int. J. Climatol.*, 22(12):1441–1453, 2002.
- 1068 S. Kotlarski, K. Keuler, O.B. Christensen, A. Colette, M. Déqué, A. Gobiet, K. Goergen,
1069 D. Jacob, D. Lüthi, E. van Meijgaard, G. Nikulin, C. Schär, C. Teichmann, R. Vautard,
1070 K. Warrach-Sagi, and V. Wulfmeyer. Regional climate modelling on European scales: A
1071 joint standard evaluation of the EURO-CORDEX RCM ensemble. *Geosci. Model. Dev.*
1072 *Discuss.*, 7:217–293, 2014.
- 1073 O. Lhotka and J. Kyseľ. Spatial and temporal characteristics of heat waves over Central
1074 Europe in an ensemble of regional climate model simulations. *Clim. Dynam.*, 45:2351–2366,
1075 2015.
- 1076 H. Li, J. Sheffield, and E.F. Wood. Bias correction of monthly precipitation and tempera-
1077 ture fields from Intergovernmental Panel on Climate Change AR4 models using equidistant
1078 quantile matching. *J. Geophys. Res.*, 115:D10101, 2010.
- 1079 P. Lorenz and D. Jacob. Validation of temperature trends in the ENSEMBLES regional
1080 climate model runs driven by ERA40. *Clim. Res.*, 44:167–177, 2010.
- 1081 D. Maraun. Bias correction, quantile mapping and downscaling: Revisiting the inflation issue.
1082 *J. Climate*, 26:2137–2143, 2013.
- 1083 D. Maraun. Bias correcting climate change simulations - a critical review. *Curr. Clim. Change*
1084 *Rep.*, 2(4):211–220, 2016. doi: 10.1007/s40641-016-0050-x.
- 1085 D. Maraun and M. Widmann. *Statistical Downscaling and Bias Correction for Climate Re-*
1086 *search*. Cambridge University Press, 2018.
- 1087 D. Maraun, F. Wetterhall, A. M. Ireson, R. E. Chandler, E. J. Kendon, M. Widmann,
1088 S. Brienen, H. W. Rust, T. Sauter, M. Themeßl, V. K. C. Venema, K. P. Chun, C. M.
1089 Goodess, R. G. Jones, C. Onof, M. Vrac, and I. Thiele-Eich. Precipitation downscaling
1090 under climate change. Recent developments to bridge the gap between dynamical models
1091 and the end user. *Rev. Geophys.*, 48:RG3003, 2010.
- 1092 D. Maraun, M. Widmann, J. M. Gutierrez, S. Kotlarski, R. E. Chandler, E. Hertig, J. Wibig,
1093 R. Huth, and R. A. I. Wilcke. VALUE: A framework to validate downscaling approaches
1094 for climate change studies. *Earth's Future*, 3:1–14, 2015.
- 1095 A. Martynov, R. Laprise, L. Sushama, K. Winger, L. Separovic, and B. Dugas. Reanalysis-
1096 driven climate simulation over CORDEX North America domain using the Canadian Re-
1097 gional Climate Model, version 5: model performance evaluation. *Clim. Dynam.*, 41:
1098 29733005, 2013.

- 1099 A. Moberg and P.D. Jones. Regional climate model simulations of daily maximum and mini-
1100 mum near-surface temperatures across Europe compared with observed station data 1961-
1101 1990. *Clim. Dynam.*, 23:695–715, 2004.
- 1102 R. Monjo, G. Chust, and V. Caselles. Probabilistic correction of RCM precipitation in the
1103 Basque Country (Northern Spain). *Theor. Appl. Climatol.*, 117(1-2):317–329, 2014.
- 1104 C: Obled, G. Bontron, and R. Garçon. Quantitative precipitation forecasts: a statistical
1105 adaptation of model outputs through an analogues sorting approach. *Atmos. Res.*, 63(3):
1106 303–324, 2002.
- 1107 R.A. Pielke and R.L. Wilby. Regional climate downscaling: What’s the point? *EOS*, 93(5):
1108 52–53, 2012.
- 1109 D.W. Pierce, D.R. Cayan, E.P. Maurer, J.T. Abatzoglou, and K.C. Hegewisch. Improved Bias
1110 Correction Techniques for Hydrological Simulations of Climate Change. *J. Hydrometeorol.*,
1111 16(6):2421–2442, 2015.
- 1112 R. Pongrácz, J. Bartholy, and A. Kis. Estimation of future precipitation conditions for Hun-
1113 gary with special focus on dry periods. *Időjárás*, 118(4):305–321, 2014.
- 1114 J. Räisänen and O. Räty. Projections of daily mean temperature variability in the future:
1115 cross-validation tests with ENSEMBLES regional climate simulations. *Clim. Dynam.*, 41:
1116 1553–1568, 2013.
- 1117 J. Rajczak, S. Kotlarski, and C. Schär. Does quantile mapping of simulated precipitation
1118 correct for biases in transition probabilities and spell lengths? *J. Climate*, 29:1605–1615,
1119 2016.
- 1120 O. Räty, J. Räisänen, and J.S. Ylhäisi. Evaluation of delta change and bias correction methods
1121 for future daily precipitation: intermodel cross-validation using ENSEMBLES simulations.
1122 *Clim. Dynam.*, 42(9-10):2287–2303, 2014.
- 1123 D. Raynaud, B. Hingray, I. Zin, S. Anquetin, S. Debionne, and R. Vautard. Atmospheric
1124 analogues for physically consistent scenarios of surface weather in Europe and Maghreb.
1125 *Int. J. Climatol.*, 37(4):2160–2176, 2017.
- 1126 C.W. Richardson. Stochastic simulation of daily precipitation, temperature, and solar radia-
1127 tion. *Wat. Resour. Res.*, 17(1), 1981.
- 1128 C. Rosenzweig, A. Iglesias, X. B. Yang, P. R. Epstein, and E. Chivian. Climate change and
1129 extreme weather events: Implications for food production, plant diseases, and pests. *Global*
1130 *Change and Human Health*, 2(2):90–104, 2001.
- 1131 C. Rosenzweig, W.D. Solecki, S.A. Hammer, and S. Mehrotra, editors. *Climate Change and*
1132 *Cities: First Assessment Report of the Urban Climate Change Research Network*. Cambridge
1133 University Press, 2011.
- 1134 M. Rummukainen. State-of-the-art with regional climate models. *Wiley Int. Rev. Clim.*
1135 *Change*, 1:82–96, 2010. DOI: 10.1002/wcc.8.

- 1136 D. San-Martín, R. Manzananas, S. Brands, S. Herrera, and J.M. Gutiérrez. Reassessing model
1137 uncertainty for regional projections of precipitation with an ensemble of statistical down-
1138 scaling methods. *J. Climate*, 30(1):203–223, 2017.
- 1139 C. Schär, D. Lüthi, U. Beyerle, and E. Heise. The soilprecipitation feedback: A process
1140 study with a regional climate model. *J. Climate*, 12(3):722–741, 1999. doi: 10.1175/1520-
1141 0442(1999)012;0722:TSPFAP;2.0.CO;2.
- 1142 A. Schindler, D. Maraun, and J. Luterbacher. Validation of the present day annual cycle in
1143 heavy precipitation over the British Islands simulated by 14 RCMs. *J. Geophys. Res.*, 117:
1144 D18107, 2007.
- 1145 J. Schmidli, C. M. Goodess, C. Frei, M. R. Haylock, Y. Hundecha, J. Ribalaygua, and
1146 T. Schmih. Statistical and dynamical downscaling of precipitation: An evaluation and com-
1147 parison of scenarios for the european alps. *Journal of Geophysical Research-Atmospheres*,
1148 112(D4), 2007.
- 1149 M. A. Semenov, R. J. Brooks, E. M. Barrow, and C. W. Richardson. Comparison of the wgen
1150 and lars-wg stochastic weather generators for diverse climates. *Clim. Res.*, 10(2):95–107,
1151 1998.
- 1152 J.C. Semenza, C.H. Rubin, K.H. Falter, J.D. Selanikio, W.D. Flanders, H.L. Howe, and J.L.
1153 Wilhelm. Heat-related deaths during the july 1995 heat wave in chicago. *N. Engl. J. Med.*,
1154 335(2):84–90, 1996.
- 1155 S. Seneviratne, D. Lüthi, M. Litschi, and C. Schär. Land-atmosphere coupling and climate
1156 change in Europe. *Nature*, 443:205–209, 2006.
- 1157 J. Sillmann, V.V. Kharin, X. Zhang, F.W. Zwiers, and D. Bronaugh. Climate extremes indices
1158 in the CMIP5 multimodel ensemble: Part 1. Model evaluation in the present climate. *J.*
1159 *Geophys. Res.*, 118(4):1716–1733, 2013.
- 1160 P. Soares and coauthors. Process based evaluation of the VALUE perfect predictor experiment
1161 of statistical downscaling methods. *Int. J. Climatol.*, *subm.*, 2017.
- 1162 P.M.M. Soares, R.M. Cardoso, P.M.A. Miranda, P. Viterbo, and M. Belo-Pereira. Assessment
1163 of the ENSEMBLES regional climate models in the representation of precipitation variability
1164 and extremes over Portugal. *J. Geophys. Res.*, 117(D7):D071114, 2012.
- 1165 P. Štěpánek, P. Zahradníček, A. Farda, P. Skalák, M. Trnka, J. Meitner, and K. Rajdl. Projec-
1166 tion of drought-inducing climate conditions in the czech republic according to euro-cordex
1167 models. *Clim. Res.*, 70(2-3):179–193, 2016.
- 1168 S. Stoll, H.J. Hendricks Franssen, M. Butts, and W. Kinzelbach. Analysis of the impact of cli-
1169 mate change on groundwater related hydrological fluxes: a multi-model approach including
1170 different downscaling methods. *Hydrology and Earth System Sciences*, 15(1):21–38, 2011.
- 1171 M. Turco, P. Quintana-Seguí, M.C. Llasat, S. Herrera, and J.M. Gutiérrez. Testing MOS
1172 precipitation downscaling for ENSEMBLES regional climate models over Spain. *J. Geophys.*
1173 *Res.*, 116(D18), 2011.

- 1174 M. Turco, M.C. Llasat, S. Herrera, and J.M. Gutiérrez. Bias correction and downscaling of
1175 future RCM precipitation projections using a MOS-Analog technique. *J. Geophys. Res.*,
1176 122(5):2631–2648, 2017.
- 1177 P. Vaittinada Ayar, M. Vrac, S. Bastin, J. Carreau, M. Déqué, and C. Gallardo. Intercompari-
1178 son of statistical and dynamical downscaling models under the EURO-and MED-CORDEX
1179 initiative framework: present climate evaluations. *Clim. Dynam.*, 46(3-4):1301–1329, 2016.
- 1180 E. van Meijgaard, L. H. van Ulft, W. J. van de Berg, F. C. Bosveld, B. J. J. M. van den
1181 Hurk, G. Lenderink, and A. P. Siebesma. The KNMI regional atmospheric climate model
1182 RACMO version 2.1. Technical Report 302, Royal Dutch Meteorological Institute, KNMI,
1183 Postbus 201, 3730 AE, De Bilt, The Netherlands, 2008.
- 1184 R. Vautard, A. Gobiet, D. Jacob, M. Belda, A. Colette, M. Déqué, J. Fernández, M. García-
1185 Díez, K. Goergen, I. Güttler, T. Halenka, T. Karacostas, E. Katragkou, K. Keuler, S. Kot-
1186 larski, S. Mayer, E. van Meijgaard, G. Nikulin, M. Patarcic, J. Scinocca, S. Sobolowski,
1187 M. Suklitsch, C. Teichmann, K. Warrach-Sagi, V. Wulfmeyer, and P. Yiou. The simulation
1188 of European heat waves from an ensemble of regional climate models within the EURO-
1189 CORDEX project. *Clim. Dynam.*, 41:2555–2575, 2013.
- 1190 C. Volosciuk, D. Maraun, V.A. Semenov, and W. Park. Extreme precipitation in an atmo-
1191 sphere general circulation model: impact of horizontal and vertical model resolutions. *J.*
1192 *Climate*, 28(3):1184–1205, 2015.
- 1193 C. Volosciuk, D. Maraun, M. Vrac, and M. Widmann. A combined statistical bias correction
1194 and stochastic downscaling method for precipitation. *Hydrol. Earth Syst. Sci.*, 21(3):1693–
1195 1719, 2017.
- 1196 H. von Storch. On the use of "inflation" in statistical downscaling. *J. Climate*, 12(12):3505–
1197 3506, 1999.
- 1198 M. Vrac and P. Friederichs. Multivariate-intervariable, spatial, and temporal-bias correction.
1199 *J. Climate*, 28(1):218–237, 2015.
- 1200 M. Vrac, P. Drobinski, A. Merlo, M. Herrmann, C. Lavaysee, L. Li, and S. Somot. Dynamical
1201 and statistical downscaling of the French Mediterranean climate: uncertainty assessment.
1202 *Nat. Haz. Earth Syst. Sci.*, 12:2769–2784, 2012.
- 1203 K. Warrach-Sagi, T. Schwitalla, V. Wulfmeyer, and H.-S. Bauer. Evaluation of a climate
1204 simulation in Europe based on the WRF–NOAH model system: precipitation in Germany.
1205 *Clim. Dynam.*, 41(3):755–774, 2013. doi: 10.1007/s00382-013-1727-7.
- 1206 M. Widmann and coauthors. Validation of spatial variability in downscaling results from the
1207 value perfect predictor experiment. *Int. J. Climatol.*, *subm.*, 2017.
- 1208 R.L. Wilby, T.M.L. Wigley, D. Conway, P.D. Jones, B.C. Hewitson, J. Main, and D.S. Wilks.
1209 Statistical downscaling of general circulation model output: A comparison of methods. *Wat.*
1210 *Resour. Res.*, 34(11):2995–3008, 1998.
- 1211 R.A.I. Wilcke, T. Mendlik, and A. Gobiet. Multi-variable error correction of regional climate
1212 models. *Clim. Change*, 120(4):871–887, 2013.

- 1213 D. S. Wilks and R. L. Wilby. The weather generation game: a review of stochastic weather
1214 models. *Prog. Phys. Geogr.*, 23(3):329–357, 1999.
- 1215 D.S. Wilks. Use of stochastic weathergenerators for precipitation downscaling. *Wiley Inter-*
1216 *disciplinary Reviews: Climate Change*, 1(6):898–907, 2010.
- 1217 G. Wong, D. Maraun, M. Vrac, M. Widmann, J. Eden, and T. Kent. Stochastic model output
1218 statistics for bias correcting and downscaling precipitation including extremes. *J. Climate*,
1219 27:6940–6959, 2014.
- 1220 C. Yang, R. E. Chandler, and V. S. Isham. Spatial-temporal rainfall simulation using gener-
1221 alized linear models. *Wat. Resour. Res.*, 41:W11415, 2005.
- 1222 W. Yang, J. Andréasson, L.P. Graham, J. Olsson, J. Rosberg, and F. Wetterhall. Distribution-
1223 based scaling to improve usability of regional climate model projections for hydrological
1224 climate change impacts studies. *Hydrol. Res.*, 41(3-4):211–229, 2010.
- 1225 W. Yang, M. Gardelin, J. Olsson, and T. Bosshard. Multi-variable bias correction: application
1226 of forest fire risk in present and future climate in Sweden. *Natural hazards and earth system*
1227 *sciences*, 15(9):2037–2057, 2015.
- 1228 T. Zerenner, V. Venema, P. Friederichs, and C. Simmer. Downscaling near-surface atmospheric
1229 fields with multi-objective Genetic Programming. *Env. Mod. Soft.*, 84:85–98, 2016.
- 1230 **Acknowledgements**
- 1231 VALUE has been funded as EU COST Action ES1102. Participation of M. Dubrovsk and R.
1232 Huth in VALUE was supported by the Ministry of Education, Youth, and Sports of the Czech
1233 Republic under contracts LD12029 and LD12059, respectively.

Code	Tech	ST	AC	SE	Predictors	Domain	Reference
MOS							
RaiRat-M6	S	no	no	yes	temperature	gridbox	Räisänen and Rätty (2013)
RaiRat-M7	S	no	no	yes	temperature	gridbox	Räisänen and Rätty (2013)
RaiRat-M8	S	no	no	yes	temperature	gridbox	Räisänen and Rätty (2013)
SB	S	no	no	yes	temperature	gridbox	
ISI-MIP	S/PM	no	no	yes	temperature	gridbox	Hempel et al. (2013)
DBS	PM	no	no	yes	temperature	gridbox	Yang et al. (2010, 2015)
GPQM	PM	no	no	no	temperature	gridbox	Bedia et al. (2016)
EQM	QM	no	no	no	temperature	gridbox	Bedia et al. (2016)
EQMs	QM	no	no	yes	temperature	gridbox	Bedia et al. (2016)
EQM-WT	QM/WT	no	no	no	temperature	gridbox	Bedia et al. (2016)
QMm	QM	no	no	yes	temperature	gridbox	Li et al. (2010)
QMBC-BJ-PR	QM	no	no	yes	temperature	gridbox	Pongrácz et al. (2014)
							Bartholy et al. (2015)
CDFt	QM	no	no	yes	temperature	gridbox	Vrac et al. (2012)
QM-DAP	QM	no	no	yes	temperature	gridbox	Štěpánek et al. (2016)
EQM-WIC658	QM	no	no	yes	temperature	gridbox	Wilcke et al. (2013)
RaiRat-M9	QM	no	no	yes	temperature	gridbox	Räisänen and Rätty (2013)
DBBC	QM	no	no	yes	temperature	gridbox	
DBD	QM	no	no	yes	temperature	gridbox	
MOS-REG	TF	yes	no	no	temperature	4 gridboxes	Herrera et al. (2017)
FIC02T	PM/A/TF	no	no	yes	temperature	gridbox	
PP							
FIC01T	A/TF	no	no	yes	Z1000+500	nat. > gridb.	
ANALOG-ANOM	A	no	no	yes	SLP/TD/T2/U+V+Z850	continental	Vaittinada Ayar et al. (2016)
ANALOG	A	no	no	no	SLP/T2/T850+700+500/Q850+500/Z500	national	Gutiérrez et al. (2013)
							San-Martín et al. (2017)
ANALOG-MP	A	no	no	yes	Z1000+500 > U+V600/T850	nat. > gridb.	Obled et al. (2002)
							Raynaud et al. (2017)
ANALOG-SP	A	no	no	yes	Z1000+500 > T2/T2-TD	nat. > gridb.	Obled et al. (2002)
							Raynaud et al. (2017)
MO-GP	TF	no	no	no	full standard set	gridbox	Zerrenner et al. (2016)
MLR-T	TF	no	no	no	T2/SLP/U+V10m/T+Q+U+V850+700+500	gridbox	
MLR-RAN	TF	no	no	no	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-RSN	TF	no	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-ASW	TF	yes	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-ASI	TF	no	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-AAN	TF	no	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-AAI	TF	no	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-AAW	TF	yes	no	yes	Z500/T850	gridbox	Huth (2002); Huth et al. (2015)
MLR-PCA-ZTR	TF	no	no	yes	Z850/T850/R850	continental	Hertig and Jacobbeit (2008)
ESD-EOFSLP	TF/WT	no	no	yes	SLP	continental	Benestad et al. (2015a)
ESD-EOFT2	TF/WT	no	no	yes	T2	continental	Benestad et al. (2015a)
ESD-SLP	TF/WT	no	no	yes	SLP	continental	Benestad et al. (2015a)
ESD-T2	TF/WT	no	no	yes	T2	continental	Benestad et al. (2015a)
MLR	TF	no	no	no	SLP/T2/T850+700+500/Q850+500/Z500	national	Gutiérrez et al. (2013)
MLR-WT	TF/WT	yes	no	yes	SLP/T2/T850+700+500/Q850+500/Z500	national	Gutiérrez et al. (2013)
WT-WG	WT/WG	yes	no	no	SLP	national	Gutiérrez et al. (2013)
SWG	TF/WG	yes	no	yes	SLP/T2/TD/U+V+Z850	continental	Vaittinada Ayar et al. (2016)
WG							
SS-WG	WG	yes	yes	yes	NA	NA	Keller et al. (2015, 2016)
MARFI-BASIC	WG	yes	yes	yes	NA	NA	
MARFI-TAD	WG	yes	yes	yes	NA	NA	
MARFI-M3	WG	yes	yes	yes	NA	NA	
GOMEZ-BASIC	WG	yes	yes	yes	NA	NA	
GOMEZ-TAD	WG	yes	yes	yes	NA	NA	

Table 2: Participating methods for temperature. Techniques: S: additive correction; PM: parametric quantile mapping; QM: empirical quantile mapping; A: analog method; TF: regression-like transfer function; WT: weather typing; WG: weather generator. Explicitly modelled: ST: stochastic noise, AC: autocorrelation, SE: seasonality. SLP: sea level pressure, T2: 2m-temperature, T: temperature, TD: dew point temperature, Z: geopotential height, Q: specific humidity, R: relative humidity, U,V,Z: velocities. A > indicates a two-step method. For the full VALUE standard set of predictors and further details on the methods see Gutiérrez and coauthors (2017) or <http://www.value-cost.eu/validationportal/app#!downscalingmethod>.

Code MOS	Tech	ST	AC	SE	Predictors	Domain	Reference
Ratyetal-M6	S	no	no	yes	precipitation	gridbox	Räty et al. (2014)
Ratyetal-M7	S	no	no	yes	precipitation	gridbox	Räty et al. (2014)
ISI-MIP	S/PM	no	no	yes	precipitation	gridbox	Hempel et al. (2013)
DBS	PM	no	no	yes	precipitation	gridbox	Yang et al. (2005, 2015)
Ratyetal-M9	PM	no	no	yes	precipitation	gridbox	Räty et al. (2014)
BC	PM	no	no	yes	precipitation	gridbox	Monjo et al. (2014)
GQM	PM	no	no	no	precipitation	gridbox	Bedia et al. (2016)
GPQM	PM	no	no	no	precipitation	gridbox	Bedia et al. (2016)
EQM	QM	no	no	no	precipitation	gridbox	Bedia et al. (2016)
EQMs	QM	no	no	yes	precipitation	gridbox	Bedia et al. (2016)
EQM-WT	QM/WT	no	no	no	precipitation	gridbox	Bedia et al. (2016)
QMm	QM	no	no	yes	precipitation	gridbox	Li et al. (2010)
QMBC-BJ-PR	QM	no	no	yes	precipitation	gridbox	Pongrácz et al. (2014)
							Bartholy et al. (2015)
CDFt	QM	no	no	yes	precipitation	gridbox	Vrac et al. (2012)
QM-DAP	QM	no	no	yes	precipitation	gridbox	Štěpánek et al. (2016)
EQM-WIC658	QM	no	no	yes	precipitation	gridbox	Wilcke et al. (2013)
Ratyetal-M8	QM	no	no	yes	precipitation	gridbox	Räty et al. (2014)
MOS-AN	A	no	no	no	precipitation	gridbox	Turco et al. (2011, 2017)
MOS-GLM	TF	yes	no	no	precipitation	4 gridboxes	Herrera et al. (2017)
VGLMGAMMA	TF/WG	yes	no	yes	precipitation	gridbox	Wong et al. (2014)
FIC02P	PM/A/TF	no	no	yes	precipitation	gridbox	
FIC04P	PM/A/TF	no	no	yes	precipitation	gridbox	
PP							
FIC01P	A/TF	no	no	yes	Z1000+500	nat. > gridb.	
FIC03P	A/TF	no	no	yes	U+V10m/U+V500/R850+700 > R850/Q700	nat. > gridb.	
ANALOG-ANOM	A	no	no	yes	SLP/TD/T2/U+V+Z850	continental	Vaittinada Ayar et al. (2016)
ANALOG	A	no	no	no	SLP/T2/T850+700+500/Q850+500/Z500	national	Gutiérrez et al. (2013) San-Martín et al. (2017)
ANALOG-MP	A	no	no	yes	Z1000+500 > U+V600/T850	nat. > gridb.	Obled et al. (2002)
ANALOG-SP	A	no	no	yes	Z1000+500 > T2/T2-TD	nat. > gridb.	Raynaud et al. (2017) Obled et al. (2002)
MO-GP	TF	no	no	no	full standard set	gridbox	Raynaud et al. (2017) Zerrenner et al. (2016)
GLM-P	TF	yes ³	no	no	Z500/T850	gridbox	
MLR-RAN	TF	no	no	no	Z500/T850	gridbox	
MLR-RSN	TF	no	no	yes	Z500/T850	gridbox	
MLR-ASW	TF	yes	no	yes	Z500/T850	gridbox	
MLR-ASI	TF	no	no	yes	Z500/T850	gridbox	
GLM-det	TF	no	no	no	SLP/T2/T850+700+500/Q850+500/Z500	national	San-Martín et al. (2017)
GLM	TF	yes	no	no	SLP/T2/T850+700+500/Q850+500/Z500	national	San-Martín et al. (2017)
GLM-WT	TF/WT	yes	no	yes	SLP/T2/T850+700+500/Q850+500/Z500 (WT: only SLP)	national	San-Martín et al. (2017)
WT-WG	WT/WG	yes	no	no	SLP	national	San-Martín et al. (2017)
SWG	TF/WG	yes	no	yes	SLP/T2/TD/U+V+Z850	continental	Vaittinada Ayar et al. (2016)
WG							
SS-WG	WG	yes	yes	yes	NA	NA	Keller et al. (2015, 2016)
MARFI-BASIC	WG	yes	yes	yes	NA	NA	
MARFI-TAD	WG	yes	yes	yes	NA	NA	
MARFI-M3	WG	yes	yes	yes	NA	NA	
GOMEZ-BASIC	WG	yes	yes	yes	NA	NA	
GOMEZ-TAD	WG	yes	yes	yes	NA	NA	

Table 3: Participating methods for precipitation. Techniques: S: scaling; PM: parametric quantile mapping; QM: empirical quantile mapping; A: analog method; TF: regression-like transfer function; WT: weather typing; WG: weather generator. Explicitly modelled: ST: stochastic noise, AC: autocorrelation, SE: seasonality. SLP: sea level pressure, T2: 2m-temperature, T: temperature, TD: dew point temperature, Z: geopotential height, Q: specific humidity, R: relative humidity, U,V,Z: velocities. A > indicates a two-step method. Methods included for illustrative purposes are marked in grey. For the full VALUE standard set of predictors and further details on the methods see Gutiérrez and coauthors (2017) or <http://www.value-cost.eu/validationportal/app#!downscalingmethod>.

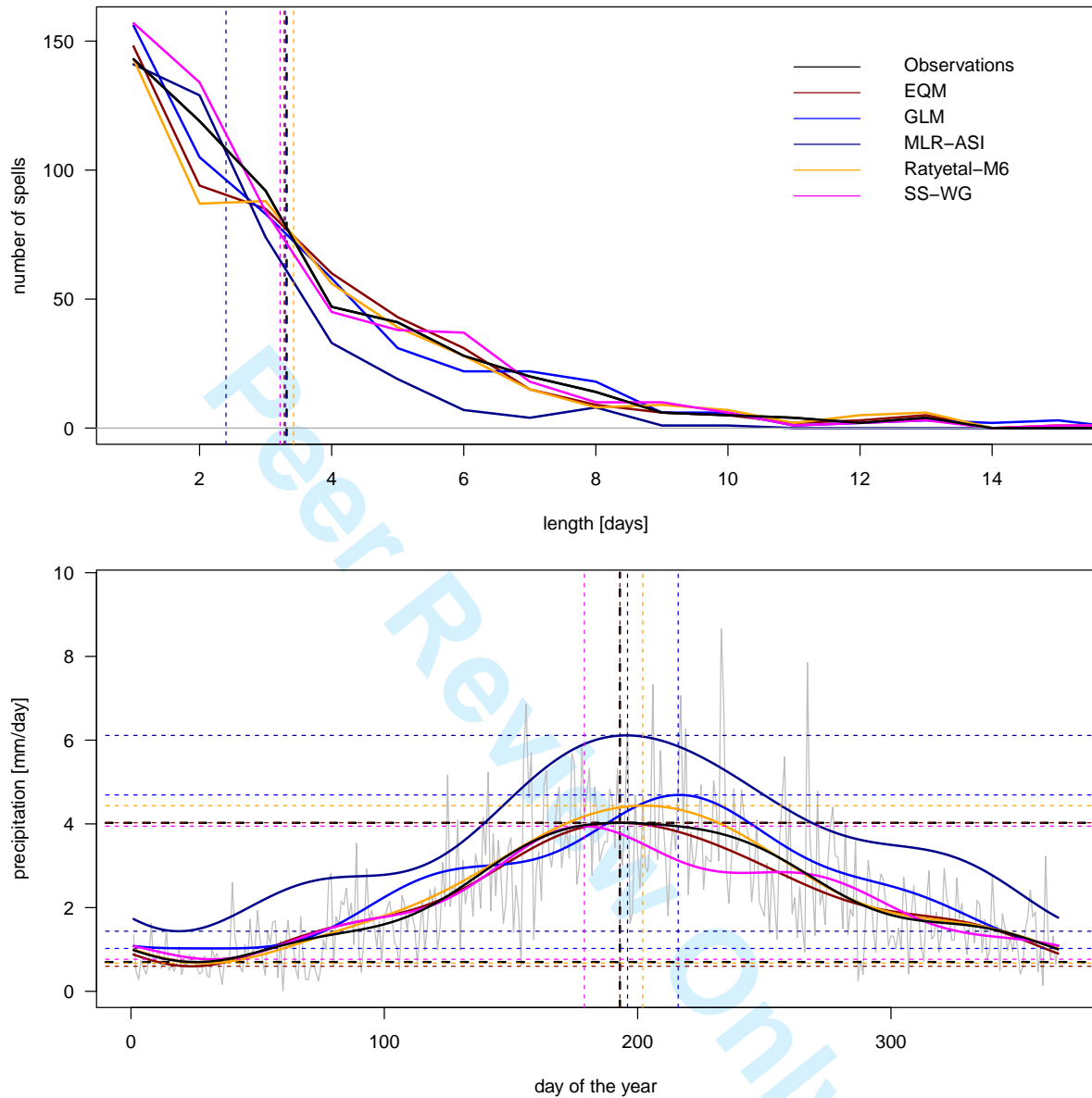


Figure 2: Illustration of selected aspects for daily precipitation, Graz, Austria. Top: dry spell length distribution. Bottom: annual cycle. Black: observations, red: EQM, orange: Ratyetal-M6, blue: MLR-SDSM, dark blue: MLR-ASI, magenta: SS-WG. Top, vertical dashed lines: mean spell length; bottom, vertical dashed lines: phase of annual cycle maximum; bottom, horizontal lines: minimum and maximum of annual cycle.

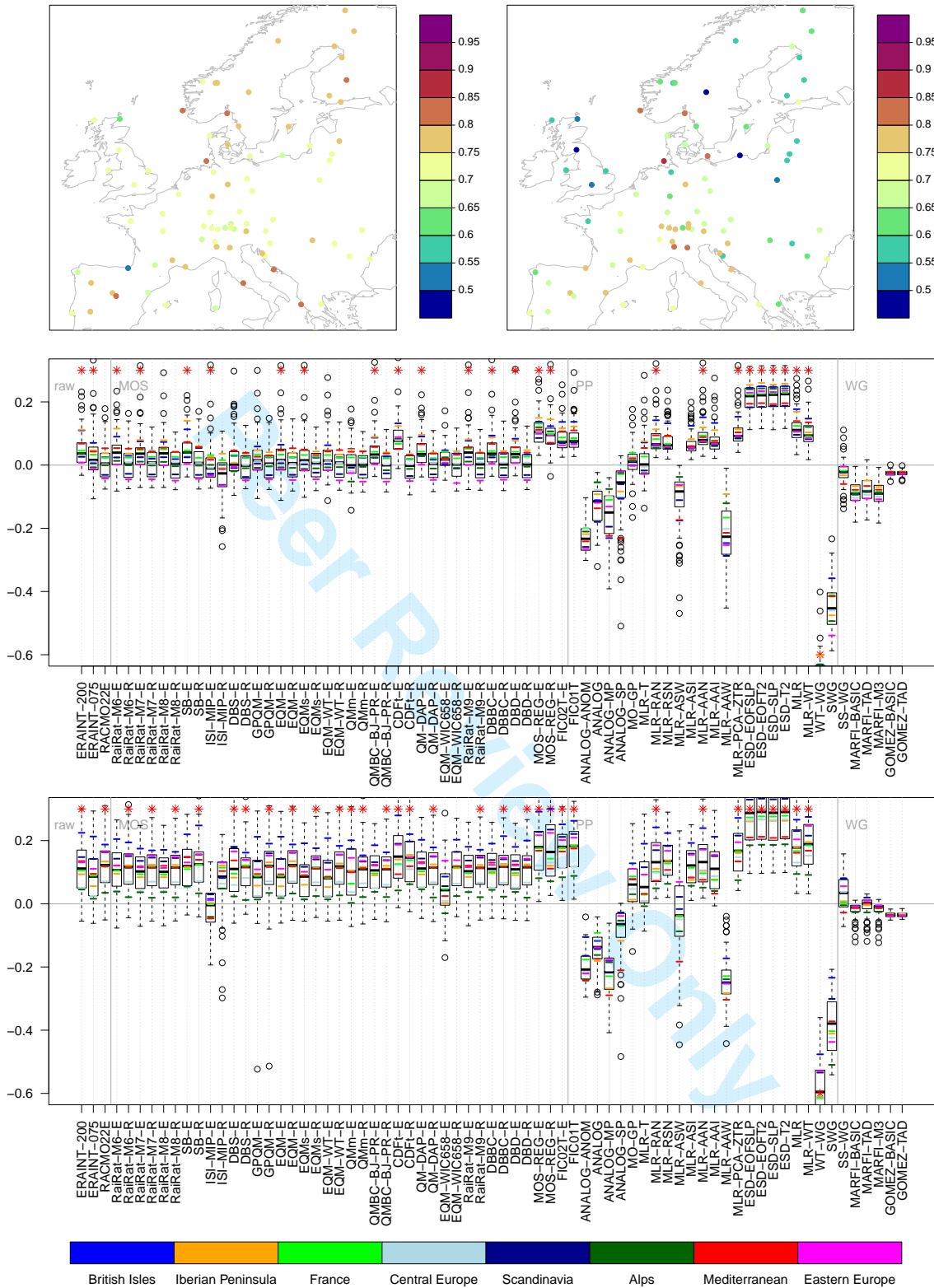


Figure 3: AC1 for summer T_{max} (left/top) and T_{min} (right/bottom). Top row: observed relationships for summer. Bottom rows: bias of the individual methods. For each method, box-whisker-plots summarise the information for all considered stations. Boxes span the 25-75% range, the whiskers the maximum value within 1.5 times the interquartile range, values outside that range are plotted individually. A red asterisk indicates that values lie outside the plotted range. The suffixes in the names of the MOS methods indicate whether a method has been driven with ERA-Interim (-E) or the RCM (-R).

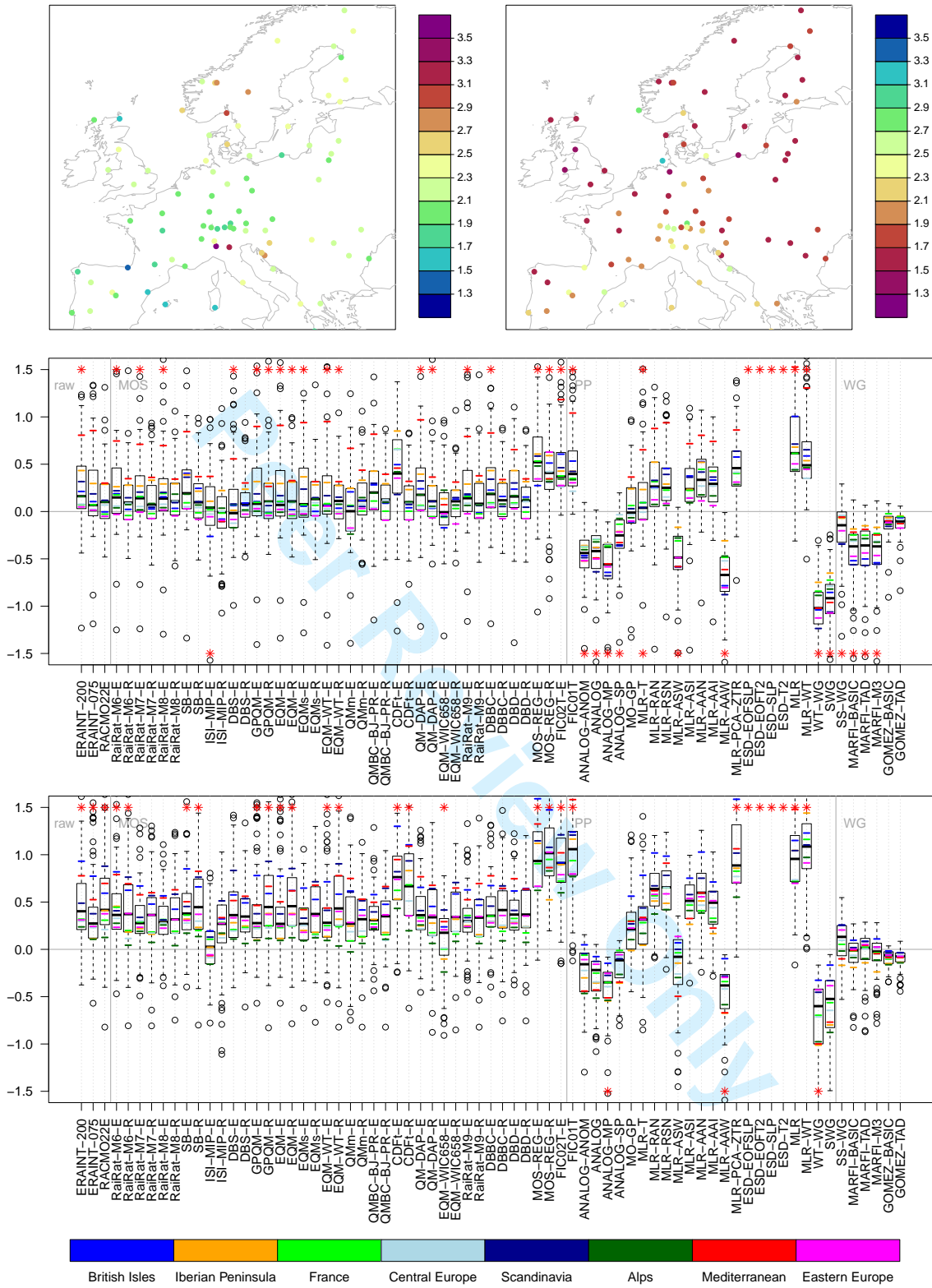


Figure 4: As Fig.3, but for summer WarmSpellMean [days] of T_{max} (top/left) and summer ColdSpellMean [days] of T_{min} (bottom/right)

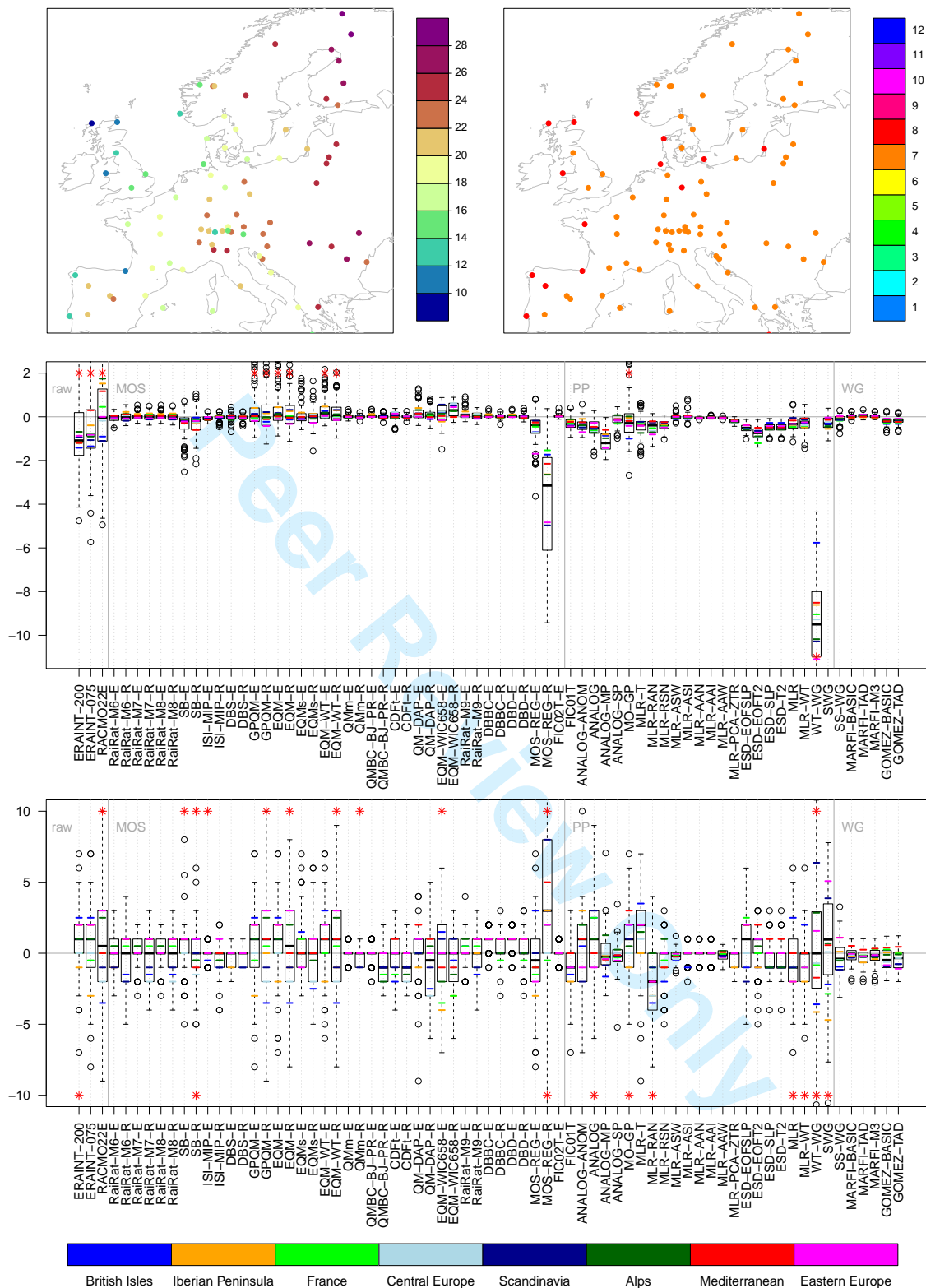


Figure 5: As Fig.3, but for the amplitude [K] (left/top) and phase [days] (right/bottom) of the annual cycle for T_{max} .

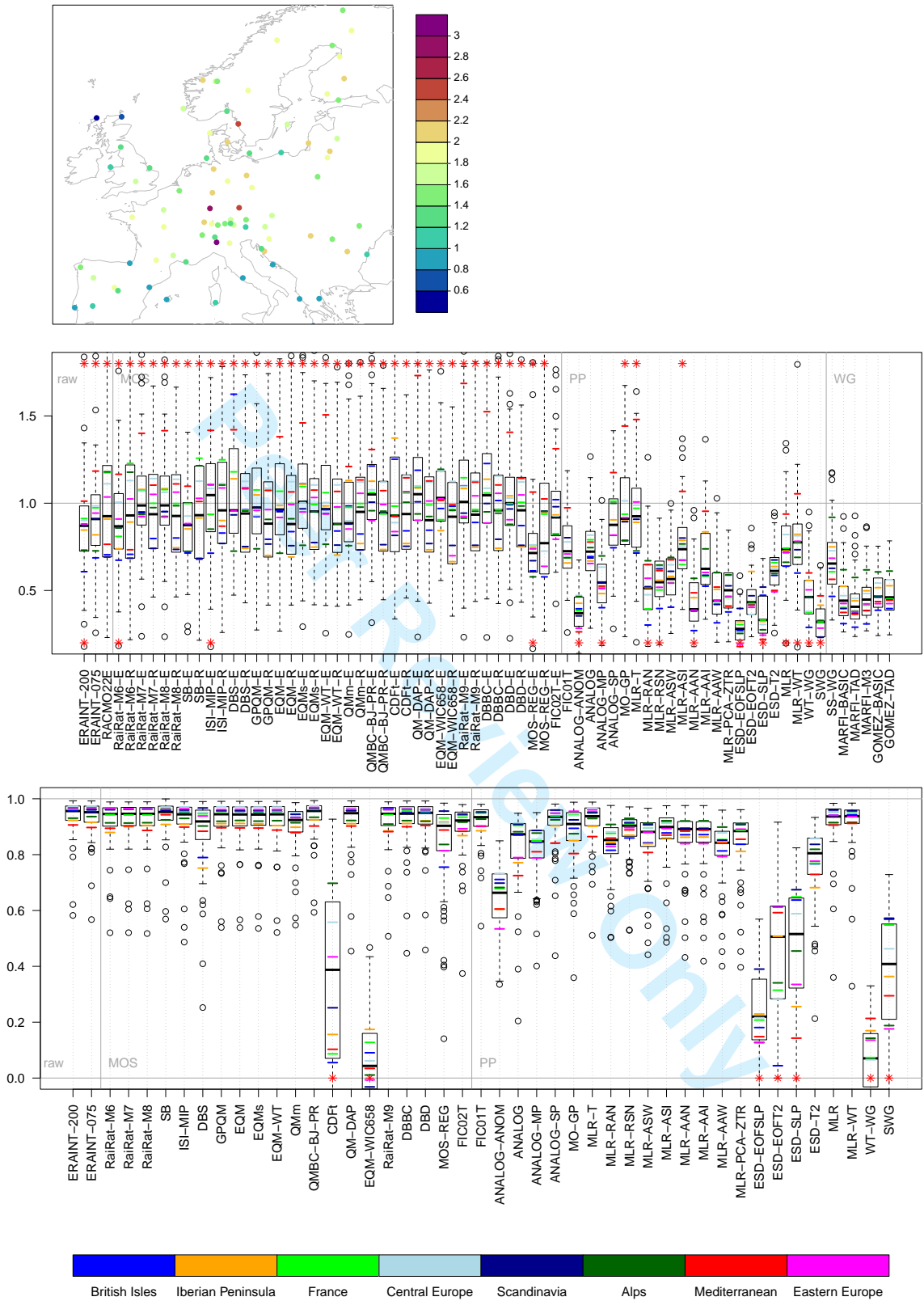


Figure 6: As Fig.3, but for summer VarY [K²] (map/top) and Cor.1Y (no map/bottom) of T_{max}.

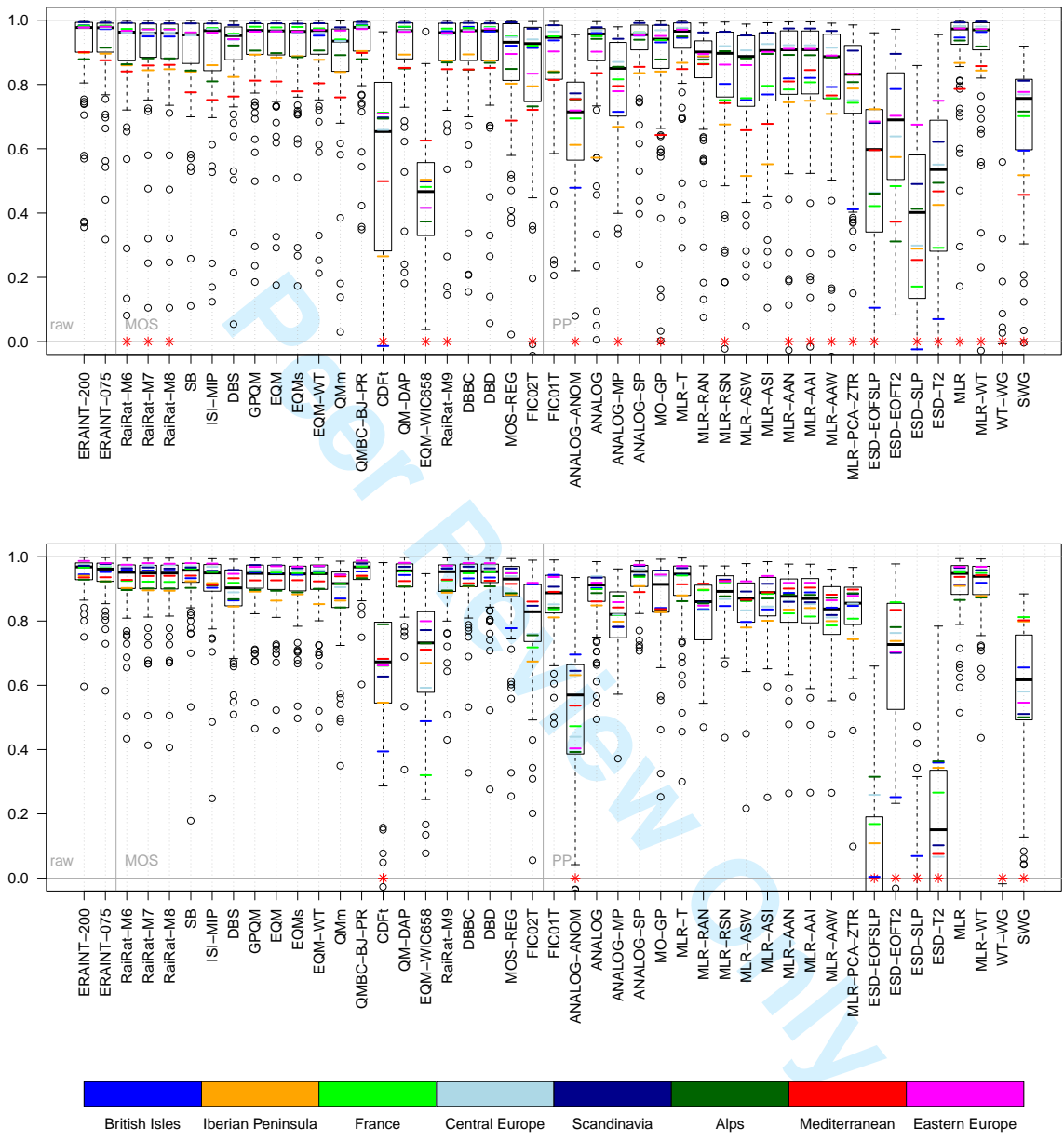


Figure 7: As Fig.3, but for Cor.7Y and T_{max} . Top: DJF; bottom: JJA.

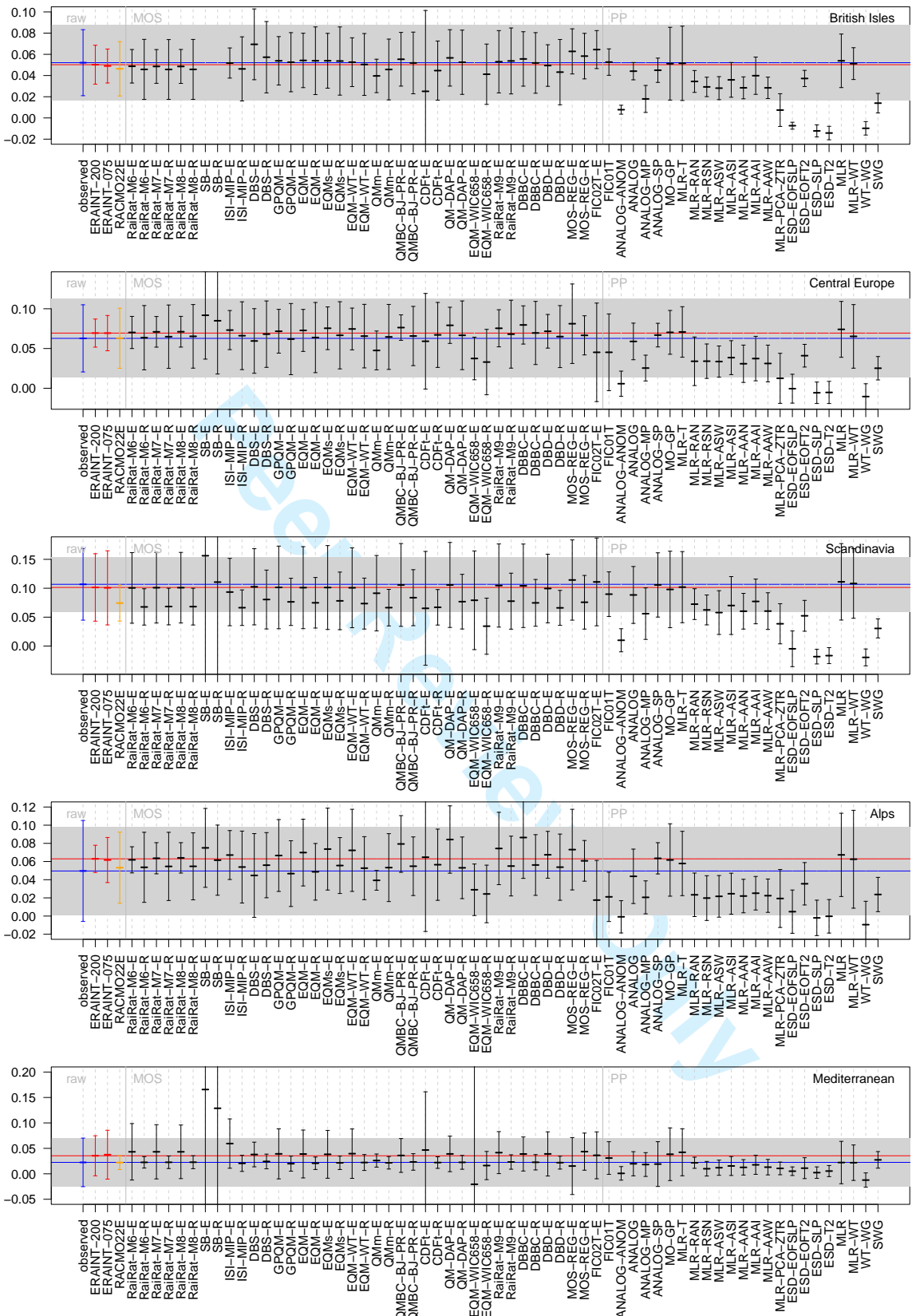


Figure 8: As Fig.3, but for the trend [K] in DJF mean T_{max} .
41

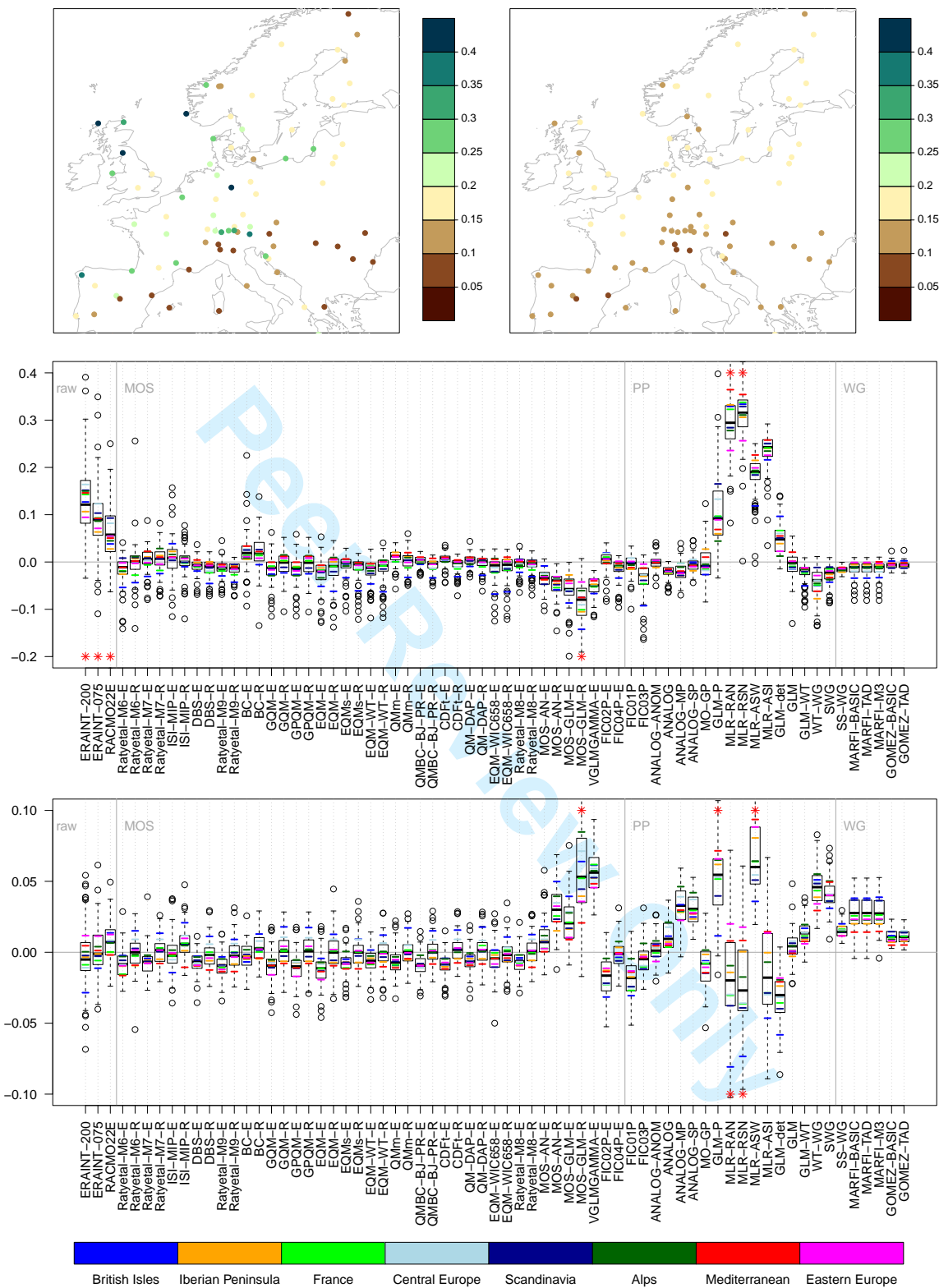


Figure 9: As Fig.3, but for winter WWProb (left/top) and DWProb (right/bottom).

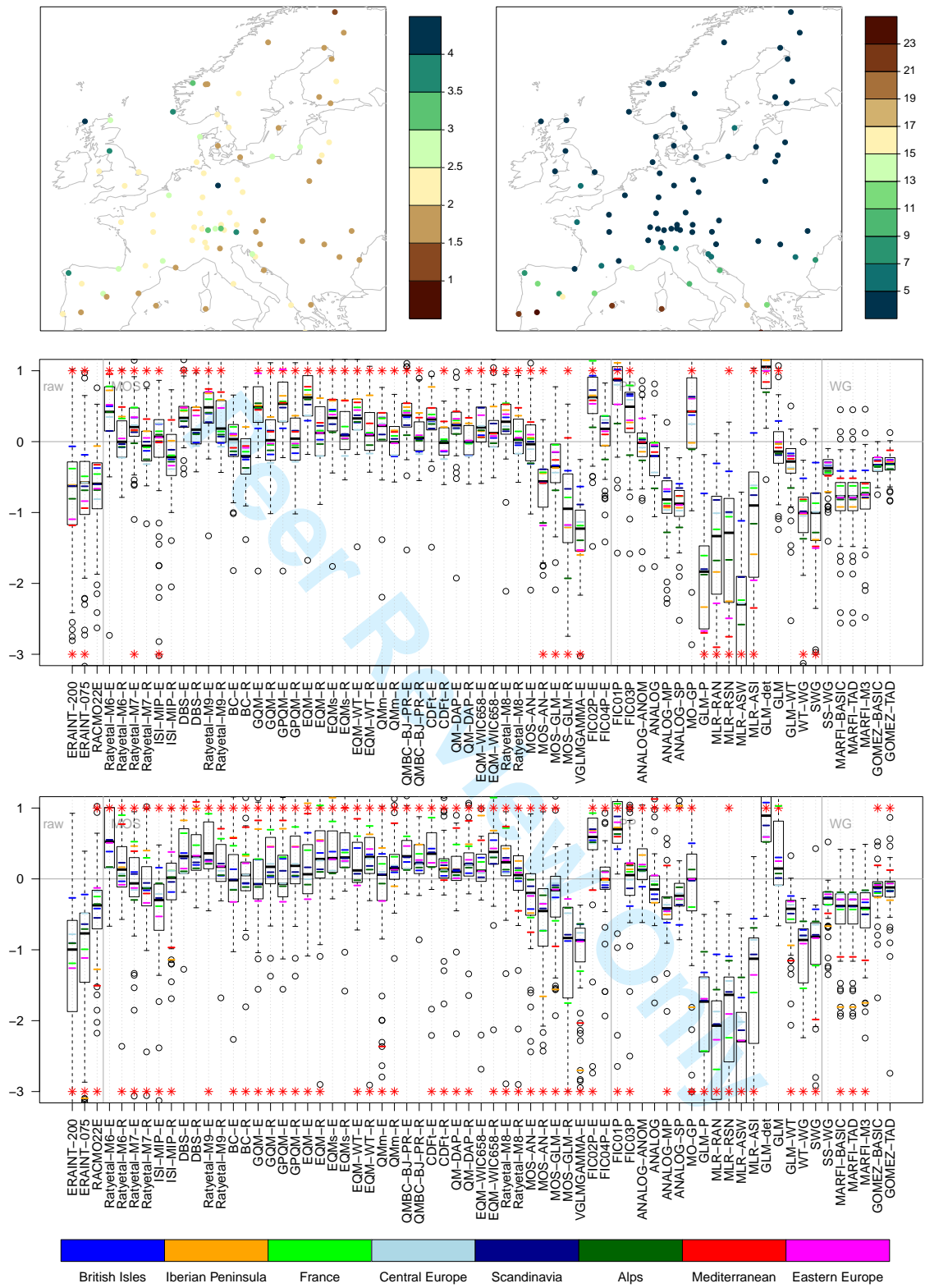


Figure 10: As Fig.3, but for winter WetSpellMean [days] (left/top) and summer DrySpellMean [days] (right/bottom)

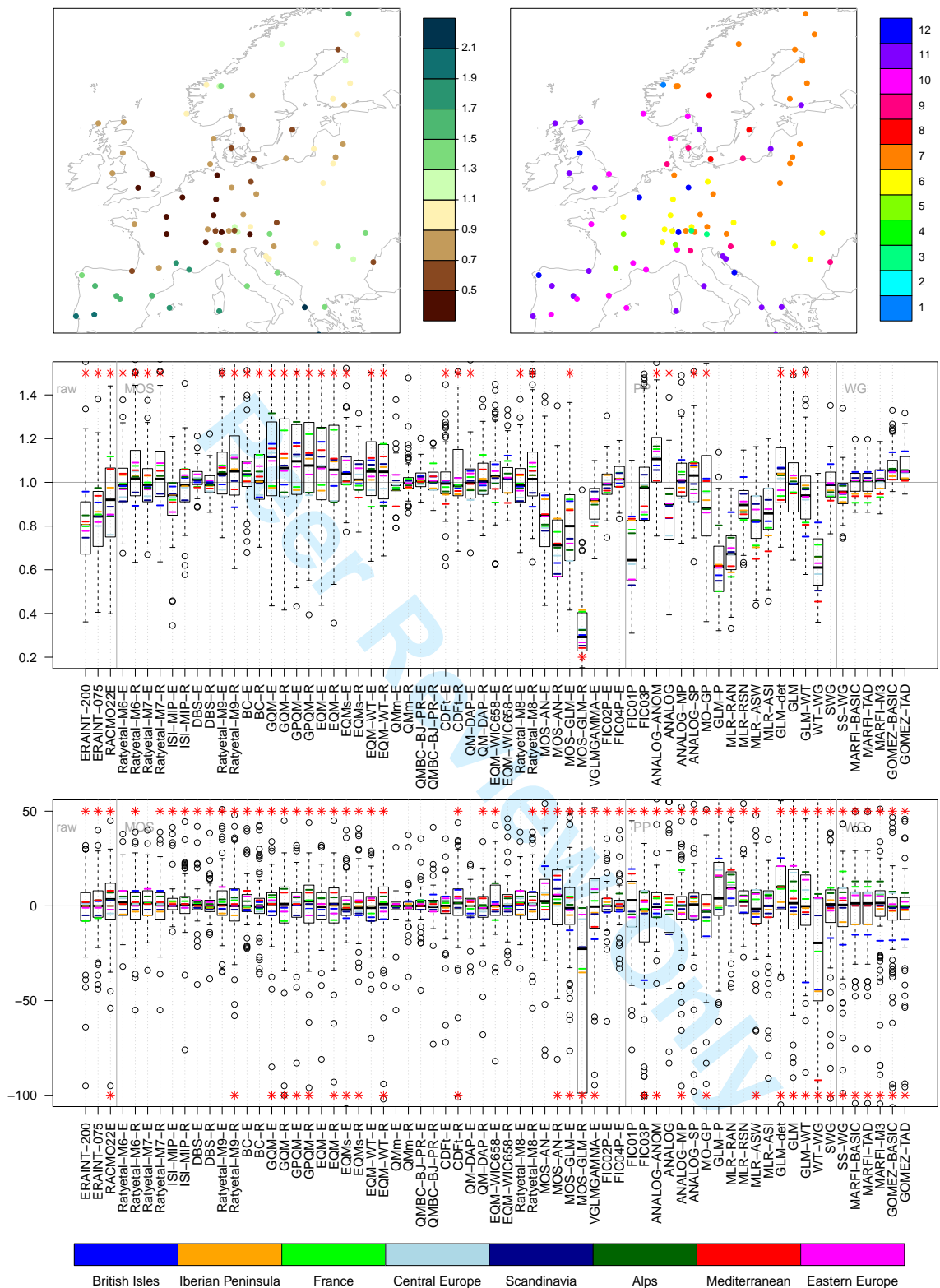


Figure 11: As Fig.3, but for the relative amplitude (left/top) and phase [days] (right/bottom) of the annual cycle of precipitation.

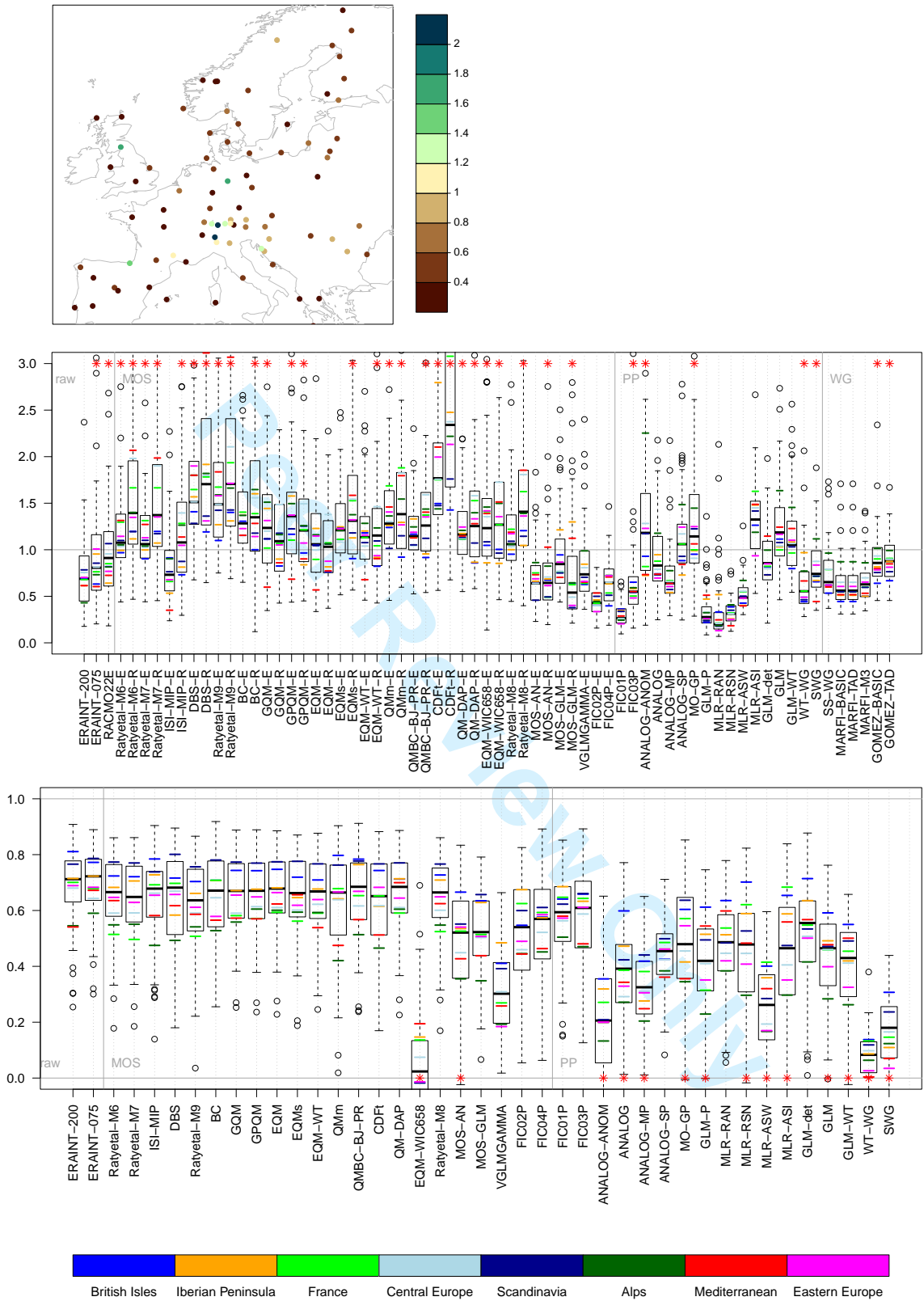


Figure 12: As Fig.3, but for summer VarY [mm²] (map/top) and Cor.1Y (no map/bottom) of precipitation.

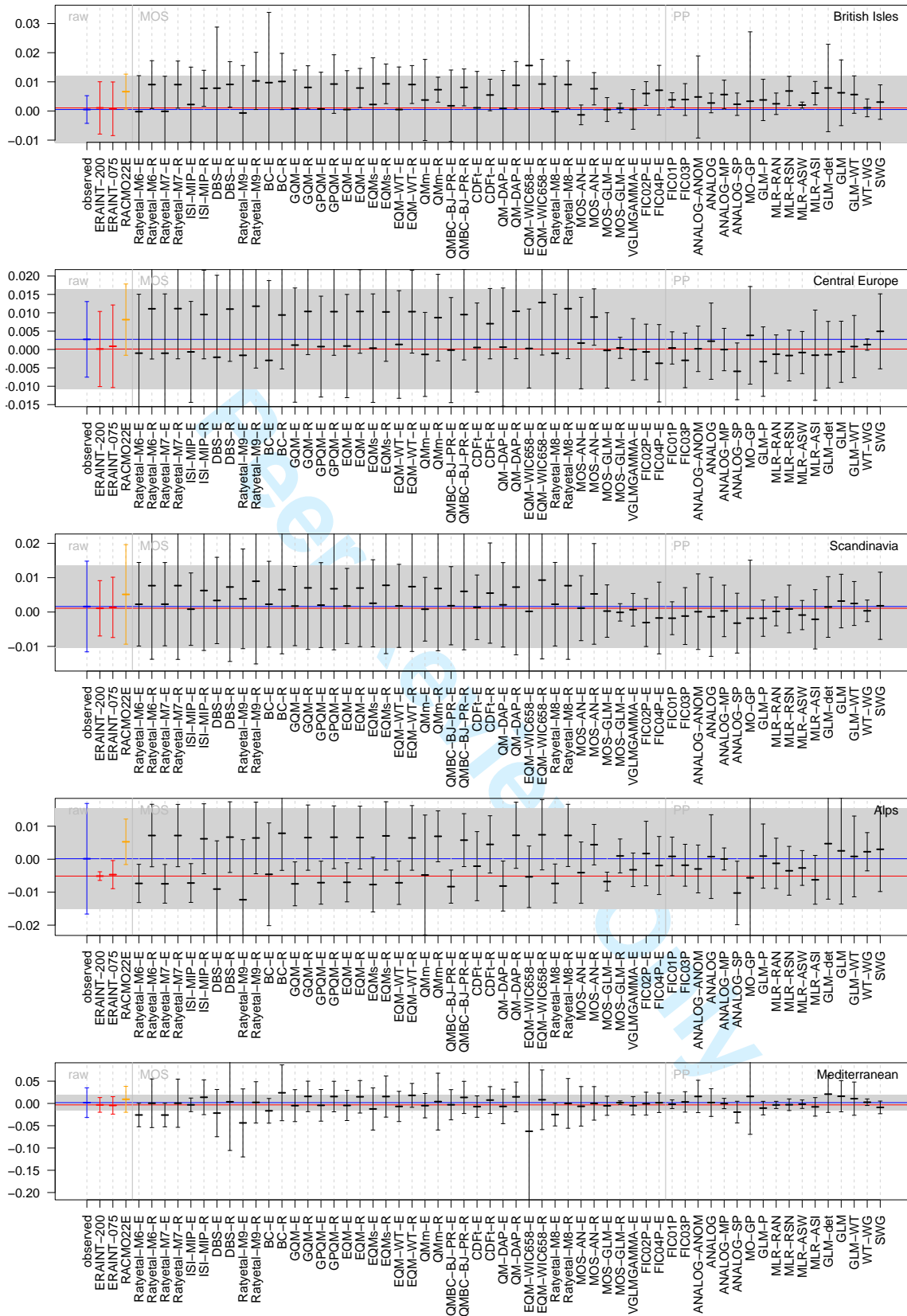


Figure 13: As Fig.3, but for the relative trend in JJA mean precipitation.

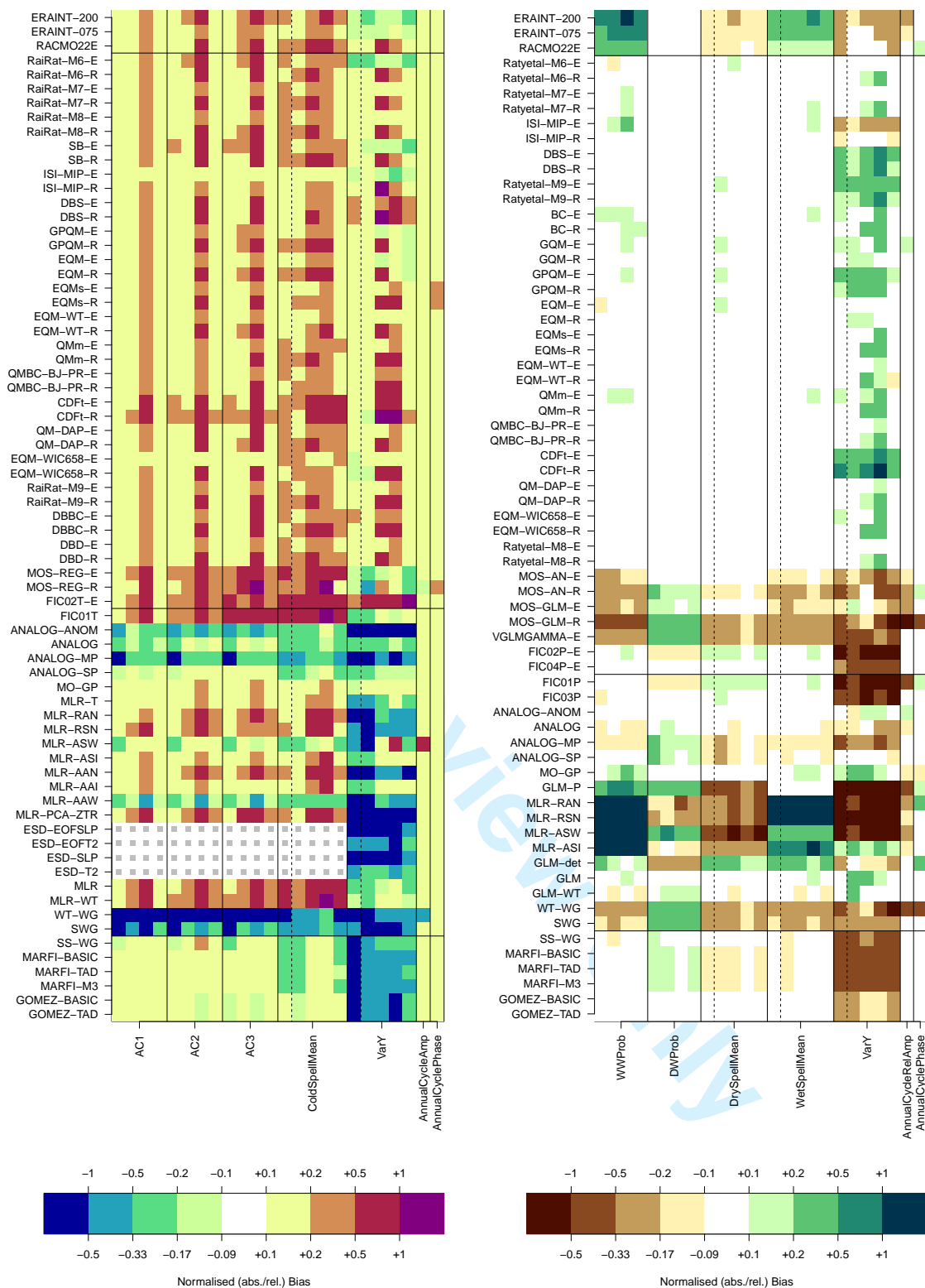


Figure 14: Performance summary. Left: T_{min} , right: precipitation. For each index either the performance for all 4 seasons is shown, or additionally the performance for the whole year (separated by a dashed line), or - in case of the seasonal cycle - only for the whole year. Grey squares indicate that no values have been calculated. For the scales used for normalisation, see Appendix.