

1 **Bias adjustment and ensemble recalibration**  
2 **methods for seasonal forecasting: A comprehensive**  
3 **intercomparison using the C3S dataset**

4 **R. Manzananas · J. M. Gutiérrez · J.**  
5 **Bhend · S. Hemri · F. J. Doblas-Reyes ·**  
6 **V. Torralba · E. Penabad · A. Brookshaw**

7  
8 Received: date / Accepted: date

9 **Abstract** This work presents a comprehensive intercomparison of different alter-  
10 natives for the calibration of seasonal forecasts, ranging from simple bias adjust-  
11 ment (BA) —e.g. quantile mapping— to more sophisticated ensemble recalibra-  
12 tion (RC) methods —e.g. non-homogeneous Gaussian regression,— which build  
13 on the temporal correspondence between the climate model and the correspond-  
14 ing observations to generate reliable predictions. To be as critical as possible, we  
15 validate the raw model and the calibrated forecasts in terms of a number of metrics  
16 which take into account different aspects of forecast quality (association, accuracy,  
17 discrimination and reliability). We focus on one-month lead forecasts of precipi-  
18 tation and temperature from four state-of-the-art seasonal forecasting systems,  
19 three of them included in the Copernicus Climate Change Service (C3S) dataset  
20 (ECMWF-SEAS5, UK Met Office-GloSea5 and Météo France-System5) for boreal  
21 winter and summer over two illustrative regions with different skill characteristics  
22 (Europe and Southeast Asia). Our results indicate that both BA and RC methods  
23 effectively correct the large raw model biases, which is of paramount importance  
24 for users, particularly when directly using the climate model outputs to feed their  
25 impact models, or when computing climate indices depending on absolute val-  
26 ues/thresholds. However, except for particular regions and/or seasons (typically  
27 with high skill), there is only marginal added value beyond this bias removal. For  
28 those cases, RC methods can outperform BA ones, mostly due to an improvement  
29 in reliability. Finally, we also show that whereas an increase in the number of mem-  
30 bers only modestly affects the results obtained from calibration, longer hindcast  
31 periods lead to improved forecast quality, particularly for RC methods.

---

R. Manzananas (✉) · J. M. Gutiérrez  
Meteorology Group, Institute of Physics of Cantabria (IFCA), CSIC-University of Cantabria  
Santander, 39005, Spain. E-mail: rmanzananas@ifca.unican.es

J. Bhend · S. Hemri  
Federal Office of Meteorology and Climatology MeteoSwiss, Switzerland

F. J. Doblas-Reyes · V. Torralba  
Barcelona Supercomputing Center (BSC), Spain

E. Penabad · A. Brookshaw  
European Centre for Medium-Range Weather Forecasts (ECMWF), UK

32 **Keywords** Seasonal forecasting, C3S, bias adjustment, ensemble recalibration,  
33 forecast quality, reliability, ensemble size, hindcast length

## 34 1 Introduction

35 The current state-of-the-art General Circulation Models (GCMs) used for sea-  
36 sonal forecasting have horizontal resolutions which are typically coarser than those  
37 needed for practical applications, and suffer from substantial systematic biases and  
38 drifts (see Doblas-Reyes et al, 2013, and references therein). This poses a risk to the  
39 use of raw model outputs in many sectors which require seasonal predictions with  
40 similar statistical properties to those observed at the regional/local scale (e.g. en-  
41 ergy, hydrology, agriculture or health). Nowadays, it is well established that some  
42 form of post-processing is needed prior to the direct use of raw model seasonal  
43 data, which constitutes a challenging problem for the development of high-quality  
44 climate services (see, e.g., Torralba et al, 2017). A number of different approaches  
45 aiming at reducing the systematic model errors have been proposed, ranging from  
46 bias adjustment (BA) and ensemble recalibration (RC) methods —both acting  
47 directly on the variable of interest,— to more complex statistical downscaling  
48 techniques building on large-scale predictors (Maraun et al, 2010). Whilst statisti-  
49 cal downscaling has been extensively analyzed in the literature in the framework of  
50 seasonal forecasting (see, e.g., Gutiérrez et al, 2005; Pavan et al, 2005; Manzan-  
51 as et al, 2018; Manzan-  
52 as and Gutiérrez, 2018; Nikulin et al, 2018), little attention  
53 has been put to-date into BA and RC methods (the focus of this work).

54 BA methods adapt the raw model outputs (e.g. predicted precipitation for a  
55 target season and a lead time) towards the corresponding observational reference  
56 to make them compatible with the local climatology. This is typically done by  
57 mapping the distribution of predicted values onto the corresponding observed one,  
58 based on a sufficiently long historical/hindcast period. These techniques, which  
59 do not use information about temporal correspondence between predictions and  
60 observations, range from simple adjustments in the mean and/or variance to more  
61 complex quantile mapping alternatives which can adjust higher order moments or  
62 even the entire distribution. Whereas the former have a long tradition in seasonal  
63 forecasting (see, e.g., Barnston, 1994; Doblas-Reyes et al, 2005), the latter have  
64 been introduced in the context of climate change projections (see, e.g., Piani et al,  
65 2010) and their application in seasonal forecasting is quite recent (see, e.g., Zhao  
66 et al, 2017; Manzan-  
67 as et al, 2018; Manzan-  
68 as and Gutiérrez, 2018). One of the main  
69 advantages of BA techniques is that they can be applied to correct daily data, even  
70 for variables which do not follow standard (e.g. Gaussian) distributions, which is  
71 often required by users (note that here we focus exclusively in the adjustment  
72 of seasonal means). For this reason, quantile mapping is rapidly becoming the  
73 method of preference by operational agencies and end-users (see, e.g., Bedia et al,  
74 2018). Nevertheless, the fact that their application is straightforward in most of  
75 the cases favors the misuse of BA techniques (Maraun et al, 2017).

76 RC methods transform the raw model outputs building on the temporal cor-  
77 respondence between the ensemble mean predictions and the corresponding ob-  
78 servations (see Sansom et al, 2016, for a comprehensive review). They range from  
79 relatively simple implementations such as climate conserving recalibration —CCR  
80 (see, e.g., Doblas-Reyes et al, 2005; Weigel et al, 2009)— or the ratio of predictable

78 components —RPC (Eade et al, 2014)— to more general ensemble model output  
79 statistics (EMOS) methods such as non-homogeneous Gaussian regression (see,  
80 e.g., Gneiting et al, 2005; Tippett and Barnston, 2008; Sansom et al, 2016). The  
81 main advantage of RC techniques for seasonal forecasting is that they are designed  
82 to produce reliable predictions. However, as opposite to BA methods, they are not  
83 suitable for the adjustment of daily data —because of the lack of model skill at  
84 this particular time-scale.— Moreover, an important constraint of RC techniques  
85 in the context of seasonal forecasting is that the underlying parameters have to be  
86 estimated using a limited amount of data and, as a consequence, they are prone to  
87 overfitting. Note that, due to the enormous computational requirements and the  
88 lack of long observational datasets required to initialize the forecasting system,  
89 state-of-the-art hindcasts typically have around 30 years of data (i.e. a sample size  
90 of 30 values for calibration).

91 Recent studies have reported some limitations for BA methods (Manzanas et al,  
92 2018) for seasonal forecasting, and even the preferable choice of RC techniques  
93 (Zhao et al, 2017). However, to the authors’ knowledge, there is no comprehensive  
94 intercomparison of BA and RC methods for this type of predictions. The main  
95 goal of this paper is therefore to fill this knowledge gap. To do this, we apply a  
96 set of state-of-the-art BA and RC methods to calibrate one-month lead seasonal  
97 predictions of temperature and precipitation from four different forecast systems.  
98 Three of these systems are included in the Copernicus Climate Change Service  
99 (C3S) seasonal service (<http://climate.copernicus.eu/seasonal-forecasts>),  
100 whereas the fourth (the ECMWF System4) is used to test the sensitivity of the  
101 results obtained to the hindcast length and the ensemble size. The raw model  
102 and calibrated predictions are validated in terms of a number of verification met-  
103 rics which take into different aspects of forecast quality (association, accuracy,  
104 discrimination and reliability).

105 The paper is organized as follows. In Section 2 we describe all the data used and  
106 introduce the different BA and RC methods applied (the implementation details  
107 are given in the Annex) and the verification metrics considered. Results obtained  
108 are presented through Section 3. The main conclusions and some interesting dis-  
109 cussion are outlined in Section 4.

## 110 **2 Data and Methods**

### 111 **2.1 Data Used**

112 In this work we focus on precipitation and temperature for boreal winter (DJF) and  
113 summer (JJA) over two illustrative regions spanning tropical and extra-tropical  
114 latitudes: Europe and Southeast Asia (EU and SA, hereafter). Note that whereas  
115 low-to-moderate skill is in general acknowledged for the former, overall good skill  
116 has been documented in the latter (see, e.g., Manzanas et al, 2014).

117 We analyze one-month lead seasonal forecasts (i.e. predictions initialized in  
118 November and May for DJF and JJA, respectively) from the C3S seasonal multi-  
119 system ensemble, which consists of three state-of-the-art models with a common  
120 hindcast period of 22 years, 1993-2014 (see Table 1). For DJF, a total of 21 sea-  
121 sons are available (starting with D1993-JF1994, which we refer to as DJF 1994).  
122 Therefore, for the sake of comparability we use a common 21-year period for both

**Table 1** Seasonal hindcasts used in this study. The last two columns show the ensemble size (members) and the period covered for each dataset.

Source	Institution	Model	Code	Members	Period
ECMWF	ECMWF	System4	System4	51	1982-2016
C3S	ECMWF	SEAS5	SEAS5	25	1993-2016
C3S	UK Met Office	GloSea5	SYSTEM12	12	1993-2015
C3S	Météo France	System5	SYSTEM5	15	1993-2014

123 DJF and JJA, and only the 12 first members of each model (minimum number  
 124 of members common across all models) are considered. Additionally, we have also  
 125 used the 51-member version of the ECMWF-System4 (Molteni et al, 2011) —  
 126 the longest to-date seasonal hindcast— for testing the sensitivity of the results  
 127 obtained to the ensemble size and the hindcast length (see Section 3.3).

128 The ERA-Interim reanalysis (Dee et al, 2011) is used as observational reference  
 129 dataset for both the calibration of the BA and RC methods and also for the  
 130 verification of all seasonal forecasts involved in this work. For consistency, all the  
 131 C3S models and ERA-Interim have been bi-linearly interpolated from their distinct  
 132 native horizontal resolutions to a common  $1^\circ$  regular grid.

## 133 2.2 Bias Adjustment and Ensemble Recalibration Methods

134 Table 2 shows the BA and RC methods intercompared in this work (see the Annex  
 135 for details on the particular implementations), which have been already used in  
 136 the context of seasonal forecasting (see the references in the third column of the  
 137 table). On the one hand, two BA methods were considered: a simple mean and  
 138 variance adjustment (MVA) and a standard implementation of empirical quantile-  
 139 quantile mapping (EQM). Note that we also considered an even simpler method  
 140 consisting of adjusting only the mean; however, the results were very similar to  
 141 those obtained for MVA and are thus not shown for brevity. Also, for EQM,  
 142 we tested the suitability of both monthly and seasonal data for the mapping,  
 143 obtaining very similar conclusions in both cases. For coherence with the rest of  
 144 methods, we only show results for the case of seasonal values. On the other hand,  
 145 four RC methods were considered: climate conserving recalibration (CCR), ratio  
 146 of predictable components (RPC) and two EMOS choices using linear regression  
 147 (LR) and non-homogeneous Gaussian regression (NGR). Note that, although more  
 148 sophisticated RC methods exist (Sansom et al, 2016), we have selected here some  
 149 standard parsimonious ones (already used in seasonal forecasting studies) which  
 150 are preferable to avoid overfitting problems.

151 All the methods considered for this work (with the exception of EQM) have  
 152 been implemented in an R-package called *calibratoR* ([http://github.com/SantanderMetGroup/  
 153 calibratoR](http://github.com/SantanderMetGroup/calibratoR)), which is publicly available as part of the *climate4R* framework  
 154 (Iturbide et al, 2018). The method EQM (as well as other BA and downscal-  
 155 ing techniques) are available in the *downscaleR* package ([http://github.com/  
 156 SantanderMetGroup/downscaleR](http://github.com/SantanderMetGroup/downscaleR)), which is also part of *climate4R*. A detailed de-  
 157 scription of each method is given in the Annex. Here, all BA and RC methods have  
 158 been applied at a gridbox level under a leave-one year-out (LOO) cross-validation  
 159 scheme (Lachenbruch and Mickey, 1968). Note that proper cross-validation is

**Table 2** Bias adjustment (BA) and ensemble recalibration (RC) methods used in this work. See the Annex for implementation details.

Approach	Method	Code	Reference(s)
BA	Mean/variance adjustment	MVA	Doblas-Reyes et al (2005), Torralba et al (2017)
BA	Empirical quantile-quantile mapping	EQM	Zhao et al (2017), Manzanas et al (2018)
RC	Climate conserving recalibration	CCR	Weigel et al (2009)
RC	Ratio of predictable components	RPC	Eade et al (2014)
RC	Linear regression	LR	Marcos et al (2018)
RC	Non-homogeneous Gaussian regression	NGR	Tippett and Barnston (2008)

160 mandatory in order to avoid artificial skill (Manzanas et al, 2017), especially when  
 161 working with small sample sizes like in this case (21 years of data).

### 162 2.3 Forecast Quality Metrics

163 The validation of seasonal predictions is a multi-faceted problem, which requires  
 164 the use of several performance metrics to analyze different aspects of forecast  
 165 quality such as association, accuracy, discrimination and reliability. Association  
 166 reflects the strength of the relationship between the forecasts and the correspond-  
 167 ing observations, which is measured here by the Pearson correlation between the  
 168 ensemble mean and the observed interannual time-series.

169 Accuracy measures the average distance between forecasts and observations.  
 170 We consider here two standard scores which are typically used to characterize  
 171 this property: the Continuous Ranked Probability Score (CRPS) and the Ranked  
 172 Probability Score (RPS). The CRPS (Hersbach, 2000) is a metric that allows to  
 173 assess the performance of probabilistic forecasts of a continuous variable based on  
 174 the integrated squared difference between the observed and the predicted cumu-  
 175 lative distribution functions (which would correspond to the mean absolute error  
 176 for deterministic forecasts). To allow for direct comparison across the different BA  
 177 and RC methods, we use the associated skill score (CRPSS), which is computed as  
 178  $1 - (CRPS_{cal}/CRPS_{ref})$ , being  $CRPS_{ref}$  the CRPS obtained for the raw model  
 179 forecasts and  $CRPS_{cal}$  the one for the calibrated predictions. Thus, values above  
 180 (below) 0, indicate that the particular calibration method improves (degrades) the  
 181 raw model prediction. The RPS (Epstein, 1969) is the discrete version of the CRPS  
 182 and measures the sum of squared differences in cumulative probability space for  
 183 a multi-category probabilistic forecast (for two-category forecasts, it would be the  
 184 Brier Score). As for the case of the CRPS, we also use here the associated skill  
 185 score (the RPSS), which is computed as  $1 - (RPS_{cal}/RPS_{ref})$ .

186 Discrimination measures the ability of the forecasts to distinguish between an  
 187 event and the corresponding non-event, which is assessed here by means of the  
 188 area under the ROC curve Kharin and Zwiers (2003) (simply referred to as ROC  
 189 hereafter). Again, we use the associated skill score (ROCSS), which is computed  
 190 as  $(ROC_{cal} - ROC_{ref})/1 - ROC_{ref}$ . This metric is recommended by the Lead  
 191 Centre for the Standardized Verification System of Long Range Forecasts and  
 192 has been used in many previous studies for the verification of seasonal forecasts  
 193 (see, e.g., Manzanas et al, 2014). Note that RPS and ROC are used here for  
 194 tercile-based probabilistic predictions. In both cases, terciles are independently  
 195 computed for the observations and the predictions, which implicitly introduces

196 a bias adjustment in the forecasts. Therefore, as opposite to CRPS, these two  
197 metrics are bias-insensitive, allowing thus to explore the added value of BA and  
198 RC method beyond the expected (by construction) model bias reduction.

199 Finally, reliability measures how closely the forecast probabilities of a certain  
200 event correspond to the observed frequency of that event (for instance a particular  
201 tercile category). Here reliability is analyzed in two different ways. On the one  
202 hand, we separate the RPS into its three components (reliability, resolution and  
203 uncertainty) following the Brier decomposition introduced in Murphy (1973). On  
204 the other hand, we use the reliability categories introduced in Weisheimer and  
205 Palmer (2014), which are based on the relative position of the best-guess reliability  
206 line and the uncertainty range around it in the reliability diagram. In particular, we  
207 use the extended classification proposed by Manzananas et al (2018), which includes  
208 the five original categories —*perfect* (green), *still very useful* (blue), *marginally*  
209 *useful* (yellow) *not useful* (orange) and *dangerously useless* (red)— plus a new one  
210 —*marginally useful +* (dark yellow).—

211 Notice also that, whereas CRPS and RPS are sensitive to changes in reliability,  
212 ROC is not. Thus, the latter also allows to assess the potential usefulness of the  
213 different calibration methods beyond the (possible) gain in reliability.

## 214 **3 Results**

### 215 3.1 Validation of Raw Model Outputs

216 As a result of their limited spatial resolution and the corresponding misrepresenta-  
217 tion of important local features (e.g. complex topography and land-sea contrasts),  
218 global models typically exhibit significant mean errors (biases) when compared  
219 with observations. Figure 1 shows bias between the one-month lead ensemble mean  
220 of the four models considered and ERA-Interim for precipitation (top) and temper-  
221 ature (bottom). As explained, the common period 1994-2014 is used and only the  
222 first 12 members are considered for all models. Important variable- and season-  
223 dependent biases are found for all models, with values over 4°C (200mm/year)  
224 for temperature (precipitation) in many locations. Although there are regional  
225 differences among models, there exists a certain common spatial pattern of bias,  
226 especially over EU (being SYSTEM5 the most dissimilar model). These systematic  
227 errors are due to the important simplifications that need to be done when build-  
228 ing the global models as a consequence of the lack of observations and knowledge,  
229 which lead to important errors in circulation, energy exchanges, etc.

230 As expected (by construction), all the BA and RC methods effectively correct  
231 the raw model biases, leading to mean errors that are smaller than 15mm/year for  
232 precipitation and 0.05°C for temperature in all cases (not shown for brevity). This  
233 is of paramount importance for users, particularly when using climate model out-  
234 puts to feed their impact models, or when computing climate indices depending on  
235 absolute values/thresholds, and proves that some form of calibration is necessary  
236 prior to the direct use of the raw predictions.

237 The temporal association between the raw ensemble mean and the correspond-  
238 ing observations is a key parameter used by the RC methods in the calibration  
239 process (see the Annex). Figure 2 shows the interannual Pearson correlation be-  
240 tween the ensemble mean of the four models considered and ERA-Interim for

241 precipitation (top) and temperature (bottom). As in Figure 1, the common period  
242 1994-2014 and the first 12 available members were considered in all cases. Only  
243 significantly positive correlations at a 90% confidence level (according to a t-test)  
244 are shown. Note that, in the following, the results found for all quality metrics are  
245 only shown for these “skillful” regions. We do this in order to avoid misinterpreta-  
246 tion of the results obtained for the RC methods, which can lead to artificial skill in  
247 regions of small (or negative) raw model correlations (see, e.g., Eade et al, 2014).  
248 Correlations are higher for temperature than for precipitation, and also higher for  
249 tropical latitudes (SA) than for extratropical ones (EU). In general, all models  
250 exhibit a similar spatial pattern of correlations, particularly for temperature.

251 We want to remark that all the results shown in Figures 1 and 2 are almost  
252 identical if all the available members are considered for each model (not shown).

### 253 3.2 Performance of BA and RC Methods

254 For brevity, we first focus in this section on the illustrative case of temperature in  
255 DJF, for which the highest added value from calibration has been found (Figures  
256 3 to 7). Then, for a comprehensive analysis we summarize the results obtained for  
257 all other cases in Figures 8 and 9.

258 Figure 3 shows the results obtained for the CRPS over EU (top) and SA (bot-  
259 tom). In particular, columns 2-3 (4-7) show the CRPSS obtained for the different  
260 BA (RC) methods, computed with respect to the CRPS for the raw outputs, which  
261 is considered as reference (column 1). Thus, values above (below) 0, shown in blue  
262 (red), indicate that the particular method improves (degrades) the raw model pre-  
263 diction. As a consequence of effectively adjusting the existing model biases (see  
264 the previous section), all methods are found to clearly improve the raw forecasts  
265 in all cases. Moreover, all methods perform similarly.

266 Figure 4 is the equivalent to Figure 3, but for the RPS (raw model outputs;  
267 column 1) and RPSS (for BA and RC methods; columns 2-7). Again, within each  
268 approach (BA or RC), all methods are found to perform very similarly. However,  
269 whereas BA methods lead in general to degraded results, RC ones provide a benefit  
270 for some particular regions (this is especially visible in SA), being this a robust  
271 feature across all models. Nevertheless, as we will show later, this benefit can not  
272 be directly generalized to other variables and/or seasons.

273 To better understand the origin of this benefit found for RC methods (as  
274 compared to BA ones), Figure 5 shows the reliability (top) and resolution (bottom)  
275 components of the RPSS shown in Figure 4. For simplicity, the results for a single  
276 model (System4) are shown; however, the same conclusions hold for the rest of  
277 models. The smaller (larger) the reliability (resolution) term is, the lower the RPS  
278 is. Therefore, the darker the color, the better in both panels. This figure proves  
279 that the improvement of RPS attained by RC methods comes from an increase in  
280 reliability (see top panel), a crucial property for the usability of seasonal forecasts.

281 Figure 6 shows the ROC (and ROCSS) for the cold and warm tercile categories  
282 (T1 and T3, top and bottom) of DJF temperature over SA. As in Figures 3 and 4,  
283 the ROC found for the raw forecasts (column 1) is considered as reference for all BA  
284 and RC methods (columns 2-7). Differently to the case of the RPSS, no added value  
285 is attained for this metric, neither for BA nor for RC methods. Moreover, results  
286 are generally degraded after calibration, particularly for RC methods —as we shall

287 see later, this can be in part explained by the short hindcast available for the C3S  
288 models.— This points out the complexity and multifaceted character of verification  
289 of seasonal forecasts, which needs to be carefully performed so that results are not  
290 misinterpreted (Doblas-Reyes et al, 2005). In particular, these results suggest that  
291 both RPS and ROC are necessary to fully assess the usefulness of multi-category  
292 probabilistic predictions.

293 Finally, we analyze how association between the predictions and observations  
294 varies with calibration. For each model (in rows), the maps in the first column of  
295 Figure 7 show the interannual Pearson correlation between ERA-Interim and the  
296 raw outputs for DJF temperature over EU (top) and SA (bottom) —this has been  
297 already shown in Figure 2.— For each of the BA and RC methods (columns 2-7),  
298 results are shown as the difference (in correlation units) with respect to the maps  
299 in column 1. As for the ROC, in general all methods are shown to degrade the  
300 correlation values attained by the raw forecasts (this is more evident for RC than  
301 for BA). It is worth to mention that this degradation in correlation is a consequence  
302 of the LOO cross-validation setting used here (all BA and RC roughly maintain  
303 the correlations exhibited by the raw outputs if cross-validation is not applied; not  
304 shown). Note the importance of this result for the potential use of BA and RC  
305 methods in operational seasonal forecasting setups.

306 In order to provide a comprehensive analysis, Figures 8 (for EU) and 9 (for  
307 SA) summarize the results obtained in terms of the different skill scores consid-  
308 ered (CRPSS, RPSS, ROCSS and correlation differences; in columns) for all cases  
309 analyzed in this work. The two variables (precipitation and temperature) and sea-  
310 sons (DJF and JJA) are shown in different rows. In all cases, results for the four  
311 available models are displayed along the x-axis. For each model, the two (four) red  
312 (black) boxplots indicate the P25-75 range for each BA (RC) method, with blue  
313 corresponding to the P10-P90 range. For EU (with a low-to-moderate skill), BA  
314 methods are in general preferable, and especially the simplest MVA. As compared  
315 to the EQM, this method is found to provide a similar adjustment of biases, whilst  
316 yielding a smaller degradation of accuracy and association measures. Differently,  
317 in SA (with high skill in some regions), RC methods yield better reliability than  
318 BA ones, although the latter (and particularly MVA) are more robust in terms of  
319 accuracy and association. Nevertheless, as we will see in the next section, this can  
320 be partially due to the short hindcast period available here.

### 321 3.3 Results' Sensitivity to Hindcast Length and Ensemble Size

322 Taking into account the limited ensemble size (12 members) and hindcast length  
323 (21 years) available for this work, we analyze the robustness of the results shown in  
324 the previous sections by assessing how the different verification metrics considered  
325 may change for larger ensemble sizes and longer hindcast periods. To do this, we  
326 use the 51-member version of the System4 (the longest hindcast to-date), and  
327 consider the period 1982-2014. For the illustrative case of DJF temperature over  
328 SA, Figure 10 shows, in different panels from top to bottom, the CRPSS, RPSS,  
329 ROCSS (only for the warm tercile category) and interannual Pearson correlation  
330 obtained for three different configurations: the 12-member ensemble for the period  
331 1994-2014 used in the previous sections (top row), a 51-member ensemble for  
332 1994-2014 (middle row) and a 51-member ensemble for 1982-2014 (bottom row).



333 Whereas a larger ensemble does not play a significant role for any of the metrics  
334 analyzed (compare top and middle rows in each panel), there is a large influence  
335 coming from the length of the hindcast period available for the RPSS and the  
336 ROCSS, and, to a lesser extent, also for correlation (compare middle and bottom  
337 rows). Note that the best results for these metrics are obtained for 1982-2014,  
338 which points out the importance of having long hindcasts for suitable calibration  
339 of seasonal forecasts. On the contrary, note also that neither the ensemble size  
340 nor the hindcast length strongly affect the results obtained for the CRPSS, which  
341 indicates that small ensembles and short hindcasts (e.g. 12 members and 21 years)  
342 are enough to robustly characterize and adjust the main systematic model errors  
343 (e.g. mean biases).

344 Additionally, Figure 11 shows the reliability categories obtained for the differ-  
345 ent configurations of the System4 considered in Figure 10. For simplicity, results  
346 are only shown for the warm tercile category (T3). Reliability is computed for each  
347 of the 20 subregions introduced in Figure 1 of Sheau et al (2017), provided there  
348 is at least a 25% of points with significantly positive interannual correlations for  
349 the ensemble mean (see Figure 2). Within each subregion, we pool together all  
350 gridboxes for both observations and predictions.

351 In agreement with the results found for the decomposition of the RPSS (Figure  
352 5), Figure 11 shows that, whereas in general BA methods do not improve (or  
353 even degrade) the reliability of the raw model outputs, RC methods tend to yield  
354 better results for particular regions. Moreover, for the case of RC methods, both  
355 ensemble size and hindcast length have an impact on reliability, being the latter  
356 the dominant factor. In particular, as compared to 1994-2014, reliability is clearly  
357 improved for the case of RC methods when considering the period 1982-2010, which  
358 suggests, again, the importance of having long hindcasts for suitable calibration.

### 359 3.4 Computational Requirements

360 Although not strictly decisive from a scientific point of view, the analysis of the  
361 computational requirements demanded by the different methods is important from  
362 a practical perspective, especially regarding their potential usability for climate  
363 services and other user-tailored applications. Figure 12 shows, for the illustrative  
364 case of DJF temperature from the System4 (12-member, 21-year version), the ex-  
365 ecution times (in minutes) required by the BA and RC methods used in this work,  
366 according to their implementation in *calibratoR*. Dark (light) gray correspond to  
367 the LOO cross-validation setting used here for EU (SA) —note that computing  
368 times drastically reduce if cross-validation is not applied; not shown.— Among  
369 the BA methods, MVA is very rapid, being therefore a suitable option for real-  
370 time applications (e.g. interactive webpages). In particular, it is much much faster  
371 than EQM, which is widely used nowadays for different sectoral tasks. Among the  
372 RC methods, CCR and RPC are computationally inexpensive choices (also po-  
373 tentially exploitable in real-time applications), with LR still providing reasonable  
374 times (less than 2 minutes for EU). Differently, the long execution times required  
375 by NGR make this method unusable for real-time operations. In the light of these  
376 results, MVA and/or CCR could be considered as benchmarking methods which  
377 provide a good compromise between performance and computational cost for the  
378 calibration of seasonal mean values.

379 **4 Discussion and Conclusions**

380 This work presents a comprehensive intercomparison of different alternatives for  
381 the calibration of seasonal forecasts, ranging from simple bias adjustment (BA)  
382 to more sophisticated ensemble recalibration (RC) methods, which build on the  
383 temporal correspondence between the climate model and the corresponding ob-  
384 servations to produce reliable forecasts. A broad set of verification metrics has  
385 been applied, accounting for different aspects of forecast quality (association, ac-  
386 curacy, discrimination and reliability). We focus precipitation and temperature  
387 from the three available C3S seasonal forecasting models (ECMWF-SEAS5, UK  
388 Met Office-GloSea5 and Météo France-System5) and validate the raw and cali-  
389 brated predictions obtained for boreal winter (DJF) and summer (JJA) over two  
390 illustrative regions with different skill characteristics (Europe and Southeast Asia).

391 Our main conclusions are the following:

- 392 1) Both approaches (BA and RC) effectively correct the large biases exhibited  
393 by raw model predictions, with the corresponding improvement in bias-sensitive  
394 metrics such as the Continuous Ranked Probability Score. This is of paramount  
395 importance for users, particularly when using climate model outputs to feed their  
396 impact models, or when computing climate indices depending on absolute val-  
397 ues/thresholds, and proves that some form of calibration is needed prior to the  
398 direct use of the raw predictions in sectoral applications.
- 399 2) For particular cases, RC methods can outperform BA ones due to an improve-  
400 ment in reliability (other aspects of forecast quality remain unaltered, or are even  
401 deteriorated). However, these situations are confined to regions and seasons with  
402 high model skill (as shown here for winter temperature in Southeast Asia).
- 403 3) As a result of the leave-one year-out cross-validation setting followed here,  
404 bias-insensitive measures are in general degraded by all calibration methods (par-  
405 ticularly by RC ones), suggesting some degree of over-fitting due to the short  
406 hindcast available. A sensitivity analysis with a longer hindcast exhibited smaller  
407 degradation, enhancing the improvement of RC results. This indicates that longer  
408 hindcast periods than those available in state-of-the-art seasonal forecasting sys-  
409 tems (e.g. C3S dataset) are needed for the robust application of RC methods.
- 410 4) Within the RC approach, all methods perform similarly, so the particular im-  
411 plementation does not play a key role. Differently, within the BA approach, the  
412 EQM method (applied here to seasonal values) is found to perform worse than  
413 the simpler MVA, particularly in terms of discrimination. Moreover, there are sig-  
414 nificant differences among distinct methods in terms of computational cost, being  
415 NGR and EQM the most time-consuming ones. This may be especially relevant for  
416 the potential usability of the different methods analyzed in real-time applications  
417 for climate services.

418  
419 In this paper we have focused on the calibration of seasonal mean values using  
420 both BA and RC methods. However, as opposite to RC, one of the potential  
421 advantages of BA methods —not explored in this work— is their suitability for  
422 daily data, which is often demanded in a variety of sectoral applications in order  
423 to run impact (crop, hydrology, etc.) models or to compute specific indices (heat  
424 waves, length of growing index, thermal comfort index, fire weather index, etc.). As  
425 a future work, we plan thus to extend the analysis presented here for BA methods  
426 to the daily scale.

427 Finally, we do not analyze here the sensitivity of the results to the observational  
 428 reference used to calibrate and validate the different methods (instead, we use a  
 429 single reference dataset, ERA-Interim). However, the results and conclusions may  
 430 be sensitive to this particular choice (especially in regions with high observational  
 431 uncertainty) so we plan to undertake a proper assessment of this factor’s impact  
 432 in a future work. Additionally, note that the choice of reference may also affect  
 433 the comparison across forecasting systems. Therefore, we do not recommend to  
 434 use the results presented here for a ranking of the different models.

## 435 **Annex: Description of BA and RC Methods**

436 All the methods described in this Annex have been applied gridbox by gridbox  
 437 considering seasonal interannual series. We use the following notation:  $y_{m,t}$  and  
 438  $y'_{m,t}$  denote the original and calibrated values of the ensemble member  $m$  at time  
 439 (season/year)  $t$ ,  $\hat{y}$  is the average of the ensemble mean ( $\bar{y}_t$ ) on all times  $t$ ,  $\hat{o}$  is  
 440 the average of the observations on all times  $t$ ,  $\sigma_f$  is the standard deviation of the  
 441 complete ensemble (pooling all member interannual time-series) and  $\sigma_o$  is the stan-  
 442 dard deviation of the observed interannual time-series. Finally,  $\rho$  is the interannual  
 443 correlation between the ensemble mean and the observational reference.

### 444 Mean (and Variance) Adjustment (MVA)

445 This is the simplest adjustment method, with a long tradition in the context of  
 446 seasonal forecasting (see, e.g., Leung et al, 1999). The ensemble mean and variance  
 447 are adjusted towards the corresponding observational ones in the following form:

$$y'_{m,t} = (y_{m,t} - \hat{y}) \frac{\sigma_o}{\sigma_f} + \hat{o} \quad (1)$$

448 A simpler version consists of correcting just the mean (MA) and has the same  
 449 formulation, but excluding the term  $\sigma_o/\sigma_f$ .

### 450 Empirical Quantile Mapping (EQM)

451 Despite quantile mapping is commonly applied to daily data, it has been re-  
 452 cently used to correct monthly precipitation from the POAMA seasonal forecasting  
 453 model (Zhao et al, 2017). Here we have considered an empirical quantile map-  
 454 ping (EQM) method participating in the VALUE downscaling intercomparison  
 455 initiative (Gutiérrez et al, 2018) which has been also recently applied to correct  
 456 seasonal precipitation forecasts (Manzanas et al, 2018; Manzanas and Gutiérrez,  
 457 2018). This method calibrates the predicted empirical probability density function  
 458 (PDF) by adjusting a number of quantiles based on the empirical observed PDF  
 459 (Déqué, 2007). In particular, here we adjust percentiles 1 to 99 and linearly in-  
 460 terpolate every two consecutive percentiles inside this range. Outside this range,  
 461 a constant extrapolation (using the correction obtained for the 1st or 99th per-  
 462 centile) is applied. This method was applied here at a ensemble-wise level; that  
 463 is, the mapping was done based on all contributing members which were pooled

464 together (all members are supposed to be statistically indistinguishable). Then,  
 465 the so-obtained unique correction factor was applied to each individual member.  
 466 Note that ensemble- and member-wise approaches have been recently reported to  
 467 provide very similar results (Manzanas et al, 2018).

#### 468 Climate Conserving Recalibration (CCR)

469 Also known as variance inflation, this method was first introduced in Doblas-Reyes  
 470 et al (2005). It modifies the predictions to have the same interannual variance as  
 471 the observational reference, while preserving their interannual correlation, and can  
 472 be expressed as:

$$y'_{m,t} = \rho \frac{\sigma_o}{std(\bar{y}_t)} \bar{y}_t + \sqrt{1 - \rho^2} \frac{\sigma_o}{\sigma_f} (y_{m,t} - \bar{y}_t) + \bar{o} \quad (2)$$

473 After Weigel et al (2009), this method has been commonly referred to as climate  
 474 conserving recalibration.

#### 475 Ratio of Predictable Components (RPC)

476 We have also considered for this work the method introduced by Eade et al (2014),  
 477 which uses the ensemble to reduce noise and adjust the forecast variance so that  
 478 the ratio of predictable components in the model and in the observations is the  
 479 same (see the paper for details). In particular, they applied the following correction  
 480 to adjust seasonal forecasts of the North Atlantic Oscillation (NAO), temperature  
 481 and pressure in the North Atlantic region:

$$y'_{m,t} = \rho \frac{\sigma_o}{std(\bar{y}_t)} (\bar{y}_t - \hat{y}) + \sqrt{1 - \rho^2} \frac{\sigma_o}{\sqrt{var(y_{m,t} - \bar{y}_t)}} (y_{m,t} - \bar{y}_t) + \hat{y} + \hat{o} \quad (3)$$

#### 482 Linear Regression Recalibration (LR)

483 This method performs a linear regression between the ensemble mean (i.e. the  
 484 time-series of  $y_t$ ) and the corresponding observations:

$$o_t = \alpha + \beta \bar{y}_t + \epsilon \quad (4)$$

485 To correct the forecast variance, the standardized anomalies are rescaled by  
 486 the standard deviation of the predictive distribution from the linear fit, so  $y'_{m,t} =$   
 487  $\alpha + \beta \bar{y}_t + \gamma_t (y_{m,t} - \bar{y}_t)$ , where

$$\gamma_t = std(\epsilon_{fit}) \sqrt{1 + 1/n + \frac{(y_t - \bar{y}_t)^2}{(n-1)var(\epsilon_{obs})}}, \quad (5)$$

488  $\epsilon_{fit}$  and  $\epsilon_{obs}$  are the residuals from the regression and the observations respectively  
 489 and the number of samples used.

491 This method (Gneiting et al, 2005) uses a constant term and the ensemble mean  
 492 signal as predictors for the calibrated forecast mean and a constant term and  
 493 the ensemble spread for the inflation (shrinkage) of the ensemble spread. The  
 494 correction has the following form:

$$y'_{m,t} = \alpha + \beta(\bar{y}_t - \hat{y}) + (y_{m,t} - \bar{y}_t)\sqrt{\gamma^2 + \delta^2 \text{var}(y_t)} \quad (6)$$

495 The parameters  $\alpha$ ,  $\beta$ ,  $\gamma$  and  $\delta$  are optimized by minimizing the ensemble CRPS.  
 496 NGR approaches have been applied in many previous works, but mostly in the  
 497 context of short-term forecasts (see, e.g., Wilks and Hamill, 2007; Thorarinsdottir  
 498 and Johnson, 2012; Feldmann et al, 2015; Scheuerer and Möller, 2015; Markus  
 499 et al, 2017). To our knowledge, only Tippett and Barnston (2008) have used it in  
 500 the context of seasonal forecasting.

501 **Acknowledgements** This work has been funded by the C3S activity on Evaluation and Quality  
 502 Control for seasonal forecasts. JMG was partially supported by the project MULTI-SDM  
 503 (CGL2015-66583-R, MINECO/FEDER). FJDR was partially funded by the H2020 EUCP  
 504 project (GA 776613).

## 505 References

- 506 Barnston AG (1994) Linear statistical short-term climate predictive skill in  
 507 the northern hemisphere. *Journal of Climate* 7(10):1513–1564, DOI 10.  
 508 1175/1520-0442(1994)007<1513:LSSTCP>2.0.CO;2, URL [http://dx.doi.org/  
 509 10.1175/1520-0442\(1994\)007<1513:LSSTCP>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(1994)007<1513:LSSTCP>2.0.CO;2)
- 510 Bedia J, Golding N, Casanueva A, Iturbide M, Buontempo C, Gutiérrez JM (2018)  
 511 Seasonal predictions of Fire Weather Index: Paving the way for their operational  
 512 applicability in Mediterranean Europe. *Climate Services* 9:101 – 110, DOI [https:  
 513 //doi.org/10.1016/j.cliser.2017.04.001](https://doi.org/10.1016/j.cliser.2017.04.001), URL [http://www.sciencedirect.com/  
 514 science/article/pii/S2405880716300826](http://www.sciencedirect.com/science/article/pii/S2405880716300826)
- 515 Dee DP, Uppala SM, Simmons AJ, Berrisford P, Poli P, Kobayashi S, Andrae U,  
 516 Balmaseda MA, Balsamo G, Bauer P, Bechtold P, Beljaars ACM, van de Berg L,  
 517 Bidlot J, Bormann N, Delsol C, Dragani R, Fuentes M, Geer AJ, Haimberger L,  
 518 Healy SB, Hersbach H, Holm EV, Isaksen L, Kallberg P, Koehler M, Matricardi  
 519 M, McNally AP, Monge-Sanz BM, Morcrette JJ, Park BK, Peubey C, de Rosnay  
 520 P, Tavolato C, Thepaut JN, Vitart F (2011) The ERA-Interim reanalysis: Con-  
 521 figuration and performance of the data assimilation system. *Quarterly Journal  
 522 of the Royal Meteorological Society* 137(656):553–597, DOI 10.1002/qj.828
- 523 Déqué M (2007) Frequency of precipitation and temperature extremes over  
 524 France in an anthropogenic scenario: Model results and statistical correction  
 525 according to observed values. *Global and Planetary Change* 57(1-2):16–26,  
 526 DOI 10.1016/j.gloplacha.2006.11.030, URL [http://www.sciencedirect.com/  
 527 science/article/pii/S0921818106002748](http://www.sciencedirect.com/science/article/pii/S0921818106002748)
- 528 Doblas-Reyes FJ, Hagedorn R, Palmer TN (2005) The rationale behind the suc-  
 529 cess of multi-model ensembles in seasonal forecasting II. Calibration and com-  
 530 bination. *Tellus A* 57(3):234–252, DOI 10.1111/j.1600-0870.2005.00104.x, URL  
 531 <http://dx.doi.org/10.1111/j.1600-0870.2005.00104.x>

532 Doblás-Reyes FJ, García-Serrano J, Lienert F, Biescas AP, Rodrigues LRL (2013)  
533 Seasonal climate predictability and forecasting: Status and prospects. Wiley  
534 Interdisciplinary Reviews: Climate Change 4(4):245–268, DOI 10.1002/wcc.217,  
535 URL <http://dx.doi.org/10.1002/wcc.217>

536 Eade R, Smith D, Scaife A, Wallace E, Dunstone N, Hermanson L, Robinson N  
537 (2014) Do seasonal-to-decadal climate predictions underestimate the predictability  
538 of the real world? Geophysical Research Letters 41(15):5620–5628, DOI  
539 10.1002/2014GL061146, URL <http://DOI.wiley.com/10.1002/2014GL061146>

540 Epstein ES (1969) A scoring system for probability forecasts of ranked  
541 categories. Journal of Applied Meteorology 8(6):985–987, DOI 10.1175/  
542 1520-0450(1969)008<0985:ASSFPF>2.0.CO;2, URL [https://doi.org/10.1175/  
543 1520-0450\(1969\)008<0985:ASSFPF>2.0.CO;2](https://doi.org/10.1175/1520-0450(1969)008<0985:ASSFPF>2.0.CO;2)

544 Feldmann K, Scheuerer M, Thorarinsdottir TL (2015) Spatial postprocessing of  
545 ensemble forecasts for temperature using non-homogeneous gaussian regression.  
546 Monthly Weather Review 143(3):955–971, DOI 10.1175/MWR-D-14-00210.1,  
547 URL <https://doi.org/10.1175/MWR-D-14-00210.1>

548 Gneiting T, Raftery AE, Westveld AH, Goldman T (2005) Calibrated probabilistic  
549 forecasting using Ensemble Model Output Statistics and minimum CRPS estimation.  
550 Monthly Weather Review 133(5):1098–1118, DOI 10.1175/MWR2904.1,  
551 URL <https://doi.org/10.1175/MWR2904.1>

552 Gutiérrez JM, Cano R, Cofiño AS, Sordo C (2005) Analysis and downscaling  
553 multi-model seasonal forecasts in Peru using self-organizing maps. Tellus A  
554 57(3):435–447, DOI 10.1111/j.1600-0870.2005.00128.x

555 Gutiérrez JM, Maraun D, Widmann M, Huth R, Hertig E, Benestad R, Roessler O,  
556 Wibig J, Wilcke R, Kotlarski S, San Martín D, Herrera S, Bedia J, Casanueva A,  
557 Manzanas R, Iturbide M, Vrac M, Dubrovsky M, Ribalaygua J, Pórtoles J, Rätty  
558 O, Räisänen J, Hingray B, Raynaud D, Casado MJ, Ramos P, Zerener T, Turco  
559 M, Bosshard T, Štěpánek P, Bartholy J, Pongracz R, Keller DE, Fischer AM,  
560 Cardoso RM, Soares PMM, Czernecki B, Pagé C (2018) An intercomparison  
561 of a large ensemble of statistical downscaling methods over Europe: Results  
562 from the VALUE perfect predictor crossvalidation experiment. International  
563 Journal of Climatology pp 1–36, DOI 10.1002/joc.5462, URL [https://rmets.  
564 onlinelibrary.wiley.com/doi/abs/10.1002/joc.5462](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/joc.5462)

565 Hersbach H (2000) Decomposition of the continuous ranked probability score for  
566 ensemble prediction systems. Weather and Forecasting 15(5):559–570, DOI 10.  
567 1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2, URL [https://doi.org/  
568 10.1175/1520-0434\(2000\)015<0559:DOTCRP>2.0.CO;2](https://doi.org/10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2)

569 Iturbide M, Bedia J, Herrera S, Baño J, Fernández J, Frías MD, Manzanas R,  
570 San-Martín D, Cimedevilla E, Cofiño AS, Gutiérrez JM (2018) climate4R: An  
571 R-based framework for climate data access, post-processing and bias correction.  
572 Under review in Environmental Modelling and Software

573 Kharin VV, Zwiers FW (2003) On the ROC score of probability forecasts.  
574 Journal of Climate 16(24):4145–4150, DOI 10.1175/1520-0442(2003)016<4145:  
575 OTRSOP>2.0.CO;2, URL [http://dx.doi.org/10.1175/1520-0442\(2003\)  
576 016<4145:OTRSOP>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(2003)016<4145:OTRSOP>2.0.CO;2)

577 Lachenbruch PA, Mickey MR (1968) Estimation of error rates in discriminant  
578 analysis. Technometrics 10(1):1–11, DOI 10.2307/1266219, URL [http://www.  
579 jstor.org/stable/1266219](http://www.jstor.org/stable/1266219)

580 Leung LR, Hamlet AF, Lettenmaier DP, Kumar A (1999) Simulations of the  
581 ENSO hydroclimate signals in the Pacific Northwest Columbia river basin.  
582 *Bulletin of the American Meteorological Society* 80(11):2313–2330, DOI  
583 10.1175/1520-0477(1999)080<2313:SOTEHS>2.0.CO;2, URL [https://doi.org/  
584 10.1175/1520-0477\(1999\)080<2313:SOTEHS>2.0.CO;2](https://doi.org/10.1175/1520-0477(1999)080<2313:SOTEHS>2.0.CO;2)

585 Manzanas R, Frías MD, Cofiño AS, Gutiérrez JM (2014) Validation of 40  
586 year multimodel seasonal precipitation forecasts: The role of ENSO on the  
587 global skill. *Journal of Geophysical Research: Atmospheres* 119(4):1708–1719,  
588 DOI 10.1002/2013JD020680, URL [http://onlinelibrary.wiley.com/doi/10.  
589 1002/2013JD020680/abstract](http://onlinelibrary.wiley.com/doi/10.1002/2013JD020680/abstract)

590 Manzanas R, Gutiérrez JM, Fernández J, van Meijgaard E, Calmanti S, Mag-  
591 ariño ME, Cofiño AS, Herrera S (2017) Dynamical and statistical downscal-  
592 ing of seasonal temperature forecasts in Europe: Added value for user appli-  
593 cations. *Climate Services* DOI 10.1016/j.cliser.2017.06.004, URL [http://www.  
594 sciencedirect.com/science/article/pii/S2405880717300067](http://www.sciencedirect.com/science/article/pii/S2405880717300067)

595 Manzanas R, Lucero A, Weisheimer A, Gutiérrez JM (2018) Can bias correction  
596 and statistical downscaling methods improve the skill of seasonal precipitation  
597 forecasts? *Climate Dynamics* 50(3):1161–1176, DOI 10.1007/s00382-017-3668-z,  
598 URL <https://link.springer.com/article/10.1007/s00382-017-3668-z>

599 Manzanas R, Gutiérrez JM (2018) Process-conditioned bias correction for  
600 seasonal forecasting: a case-study with ENSO in Peru. *Climate Dynamics*  
601 pp 1–11, DOI 10.1007/s00382-018-4226-z, URL [https://doi.org/10.1007/  
602 s00382-018-4226-z](https://doi.org/10.1007/s00382-018-4226-z)

603 Maraun D, Wetterhall F, Ireson AM, Chandler RE, Kendon EJ, Widmann M,  
604 Brien S, Rust HW, Sauter T, Themessl M, Venema VKC, Chun KP, Good-  
605 erness CM, Jones RG, Onof C, Vrac M, Thiele-Eich I (2010) Precipitation down-  
606 scaling under climate change: Recent developments to bridge the gap between  
607 dynamical models and the end user. *Reviews of Geophysics* 48(3):n/a–n/a, DOI  
608 10.1029/2009RG000314, URL <http://dx.doi.org/10.1029/2009RG000314>

609 Maraun D, Shepherd TG, Widmann M, Zappa G, Walton D, M GJ, Hagemann S,  
610 Richter I, Soares PMM, Hall A, Mearns LO (2017) Towards process-informed  
611 bias correction of climate change simulations. *Nature Climate Change* 7:764–  
612 773, DOI 10.1038/nclimate3418

613 Marcos R, Llasat MC, Quintana-Seguí P, Turco M (2018) Use of bias correc-  
614 tion techniques to improve seasonal forecasts for reservoirs: A case-study in  
615 northwestern Mediterranean. *Science of The Total Environment* 610-611:64  
616 – 74, DOI <https://doi.org/10.1016/j.scitotenv.2017.08.010>, URL [http://www.  
617 sciencedirect.com/science/article/pii/S0048969717320089](http://www.sciencedirect.com/science/article/pii/S0048969717320089)

618 Markus D, J MG, W MJ, Achim Z (2017) Spatial ensemble postprocessing with  
619 standardized anomalies. *Quarterly Journal of the Royal Meteorological Society*  
620 143(703):909–916, DOI 10.1002/qj.2975, URL [https://rmets.onlinelibrary.  
621 wiley.com/doi/abs/10.1002/qj.2975](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.2975)

622 Molteni F, Stockdale T, Balsameda M, Balsamo G, Buizza R, Ferranti L,  
623 Magnusson L, Mogensen K, Palmer T, Vitart F (2011) The new ECMWF  
624 seasonal forecast system (System 4). European Centre for Medium-Range  
625 Weather Forecasts, URL [http://climate.ncas.ac.uk/people/allan/Fire\\_  
626 Risk\\_Insurance\\_Papers/Moltini%20etal%202011.pdf](http://climate.ncas.ac.uk/people/allan/Fire_Risk_Insurance_Papers/Moltini%20etal%202011.pdf)

627 Murphy AH (1973) A new vector partition of the probability score. *Journal of*  
628 *Applied Meteorology* 12(4):595–600, DOI 10.1175/1520-0450(1973)012<0595:

629 ANVPOT>2.0.CO;2, URL [https://doi.org/10.1175/1520-0450\(1973](https://doi.org/10.1175/1520-0450(1973)  
630 [012<0595:ANVPOT>2.0.CO;2](https://doi.org/10.1175/1520-0450(1973)

631 Nikulin G, Asharafb S, Magariño ME, Calmanti S, Cardoso RM, Bhend J,  
632 Fernández J, Frías MD, Fröhlichb K, Frühb B, Herrera S, Manzanar R, Gutiérrez  
633 JM, Hanssona U, Kolaxa M, Liniger M, Soares PMM, Spirig C, Tome R, Wysera  
634 K (2018) Dynamical and statistical downscaling of a global seasonal hindcast  
635 in eastern Africa. *Climate Services* 9:72 – 85, DOI [https://doi.org/10.1016/](https://doi.org/10.1016/j.cliser.2017.11.003)  
636 [j.cliser.2017.11.003](https://doi.org/10.1016/j.cliser.2017.11.003), URL [http://www.sciencedirect.com/science/article/](http://www.sciencedirect.com/science/article/pii/S2405880717300055)  
637 [pii/S2405880717300055](http://www.sciencedirect.com/science/article/pii/S2405880717300055)

638 Pavan V, Marchesi S, Morgillo A, Cacciamani C, Doblas-Reyes FJ (2005) Down-  
639 scaling of DEMETER winter seasonal hindcasts over Northern Italy. *Tellus A*  
640 57(3):424–434, DOI [10.1111/j.1600-0870.2005.00111.x](https://doi.org/10.1111/j.1600-0870.2005.00111.x)

641 Piani C, Haerter JO, Coppola E (2010) Statistical bias correction for daily precip-  
642 itation in regional climate models over Europe. *Theoretical and Applied Clima-*  
643 *tology* 99(1-2):187–192, DOI [10.1007/s00704-009-0134-9](https://doi.org/10.1007/s00704-009-0134-9), URL [http://dx.doi.](http://dx.doi.org/10.1007/s00704-009-0134-9)  
644 [org/10.1007/s00704-009-0134-9](http://dx.doi.org/10.1007/s00704-009-0134-9)

645 Sansom PG, Ferro CAT, Stephenson DB, Goddard L, Mason SJ (2016) Best prac-  
646 tices for postprocessing ensemble climate forecasts. Part I: Selecting appropri-  
647 ate recalibration methods. *Journal of Climate* 29(20):7247–7264, DOI [10.1175/](https://doi.org/10.1175/JCLI-D-15-0868.1)  
648 [JCLI-D-15-0868.1](https://doi.org/10.1175/JCLI-D-15-0868.1), URL <https://doi.org/10.1175/JCLI-D-15-0868.1>

649 Scheuerer M, Möller D (2015) Probabilistic wind speed forecasting on a grid based  
650 on ensemble model output statistics. *The Annals of Applied Statistics* 9(3):1328–  
651 1349, DOI [10.1214/15-AOAS843](https://doi.org/10.1214/15-AOAS843), URL <https://doi.org/10.1214/15-AOAS843>

652 Sheau TN, Tangang F, Juneng L (2017) Bias correction of global and regional  
653 simulated daily precipitation and surface mean temperature over Southeast  
654 Asia using quantile mapping method. *Global and Planetary Change* 149:79  
655 – 90, DOI <https://doi.org/10.1016/j.gloplacha.2016.12.009>, URL [http://www.](http://www.sciencedirect.com/science/article/pii/S0921818116301266)  
656 [sciencedirect.com/science/article/pii/S0921818116301266](http://www.sciencedirect.com/science/article/pii/S0921818116301266)

657 Thorarinsdottir TL, Johnson MS (2012) Probabilistic wind gust forecast-  
658 ing using non-homogeneous gaussian regression. *Monthly Weather Review*  
659 140(3):889–897, DOI [10.1175/MWR-D-11-00075.1](https://doi.org/10.1175/MWR-D-11-00075.1), URL [https://doi.org/10.](https://doi.org/10.1175/MWR-D-11-00075.1)  
660 [1175/MWR-D-11-00075.1](https://doi.org/10.1175/MWR-D-11-00075.1)

661 Tippett MK, Barnston AG (2008) Skill of multimodel ENSO probability forecasts.  
662 *Monthly Weather Review* 136(10):3933–3946, DOI [10.1175/2008MWR2431.1](https://doi.org/10.1175/2008MWR2431.1)

663 Torralba V, Doblas-Reyes FJ, MacLeod D, Christel I, Davis M (2017) Sea-  
664 sonal climate prediction: A new source of information for the management  
665 of wind energy resources. *Journal of Applied Meteorology and Climatology*  
666 56(5):1231–1247, DOI [10.1175/JAMC-D-16-0204.1](https://doi.org/10.1175/JAMC-D-16-0204.1), URL [https://doi.org/](https://doi.org/10.1175/JAMC-D-16-0204.1)  
667 [10.1175/JAMC-D-16-0204.1](https://doi.org/10.1175/JAMC-D-16-0204.1)

668 Weigel AP, Liniger MA, Appenzeller C (2009) Seasonal ensemble forecasts: Are  
669 recalibrated single models better than multimodels? *Monthly Weather Review*  
670 137(4):1460–1479, DOI [10.1175/2008MWR2773.1](https://doi.org/10.1175/2008MWR2773.1)

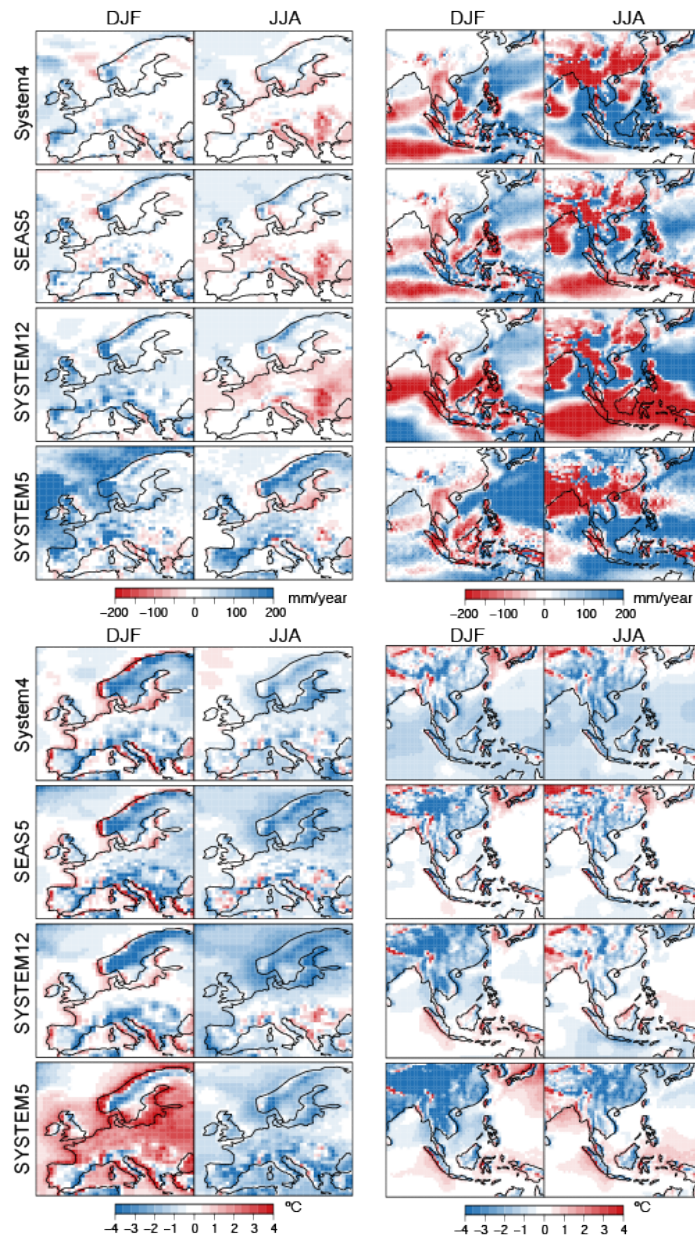
671 Weisheimer A, Palmer TN (2014) On the reliability of seasonal climate forecasts.  
672 *Journal of the Royal Society Interface* 11(96), DOI [10.1098/rsif.2013.1162](https://doi.org/10.1098/rsif.2013.1162)

673 Wilks DS, Hamill TM (2007) Comparison of Ensemble-MOS methods using  
674 GFS reforecasts. *Monthly Weather Review* 135(6):2379–2390, DOI [10.1175/](https://doi.org/10.1175/MWR3402.1)  
675 [MWR3402.1](https://doi.org/10.1175/MWR3402.1), URL <https://doi.org/10.1175/MWR3402.1>

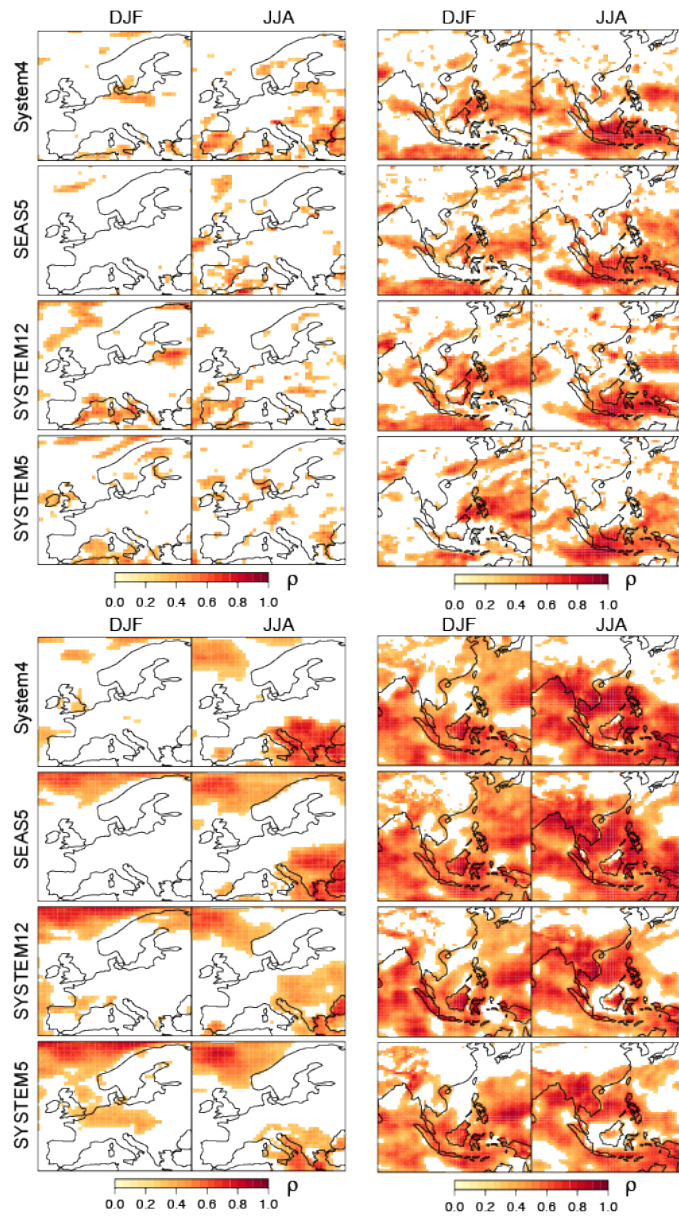
676 Zhao T, Bennett JC, Wang QJ, Schepen A, Wood AW, Robertson DE, Ramos MH  
677 (2017) How suitable is quantile mapping for postprocessing GCM precipitation



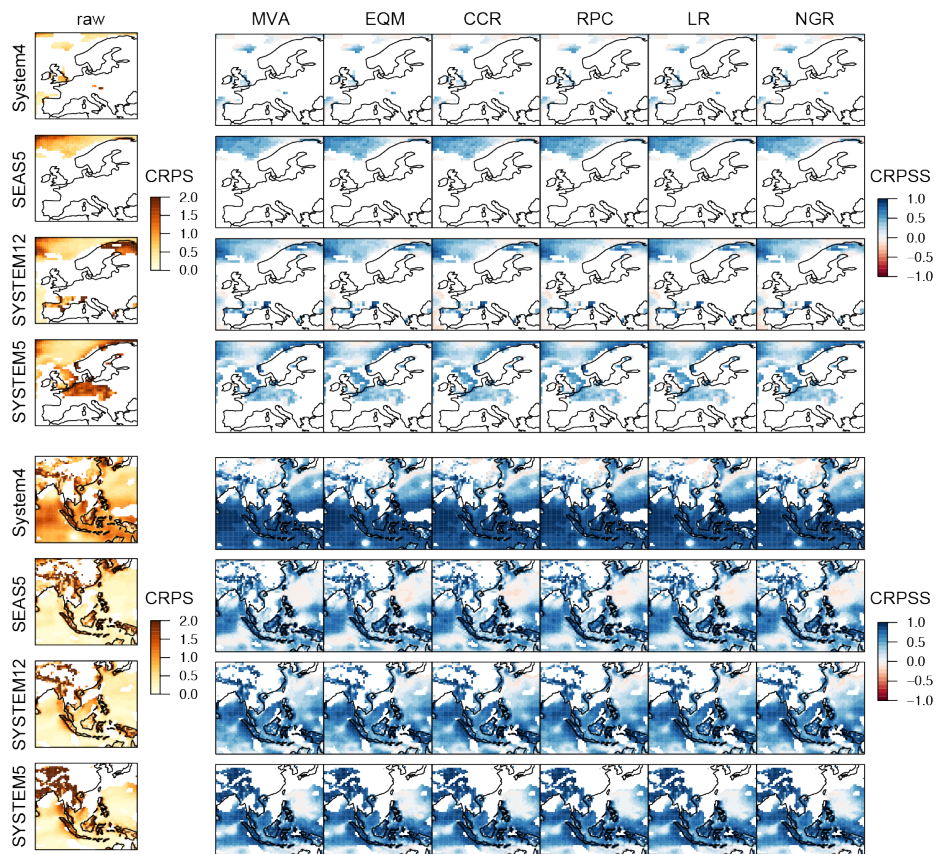
678 forecasts? *Journal of Climate* 30(9):3185–3196, DOI 10.1175/JCLI-D-16-0652.1,  
679 URL <https://doi.org/10.1175/JCLI-D-16-0652.1>



**Fig. 1** Bias between the ensemble mean of the four models of Table 1 and ERA-Interim verifying observations for precipitation (top) and temperature (bottom) over EU (left) and SA (right), in DJF and JJA. The errors are expressed as mm/year ( $^{\circ}\text{C}$ ) for the case of precipitation (temperature).



**Fig. 2** As Figure 1 but for interannual Pearson correlation. Only significant correlations (90% confidence level, according to a t-test) are shown.



**Fig. 3** CRPSS for temperature over EU (top) and SA (bottom) in DJF, as obtained from applying the different BA and RC methods of Table 2 (columns 2-7) to the four models of Table 1 (in rows). In all cases, the CRPS obtained for the raw outputs (column 1) is considered as reference.

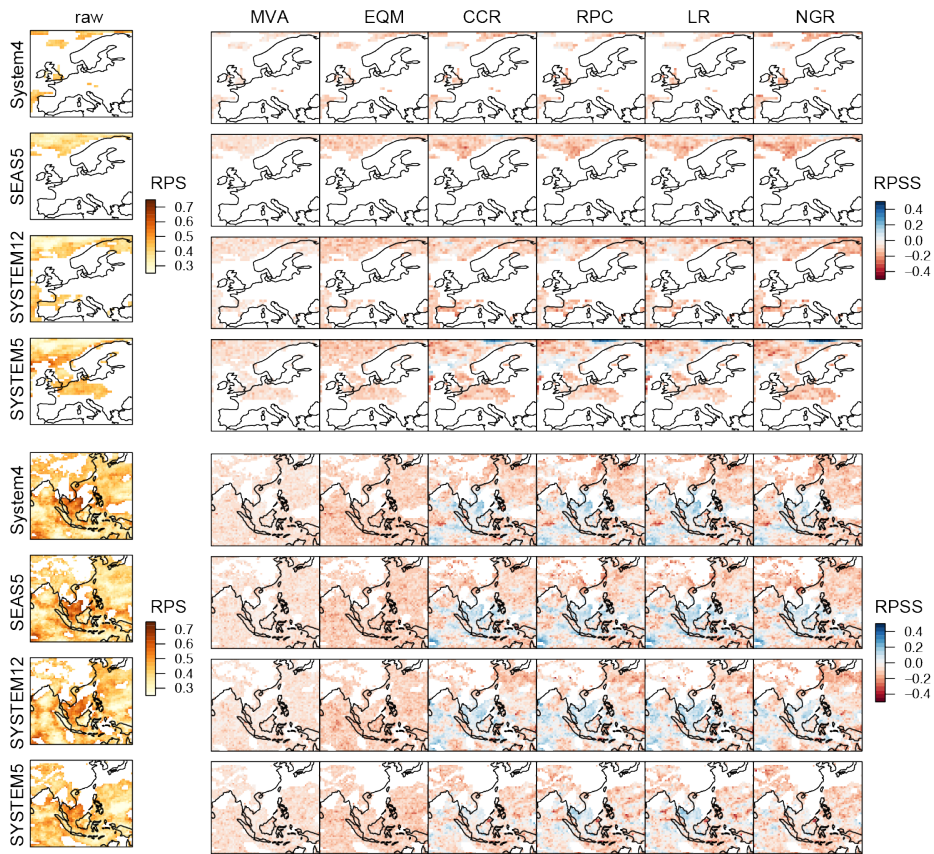


Fig. 4 As Figure 3, but for the RPSS.

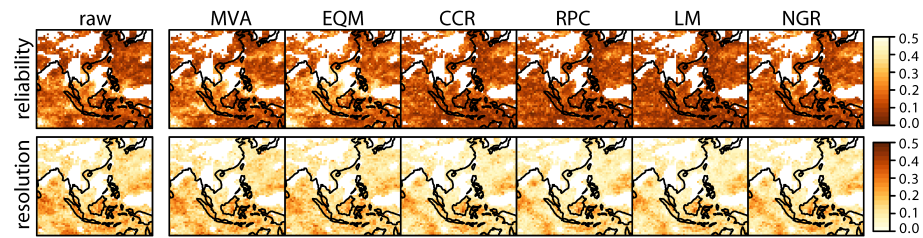
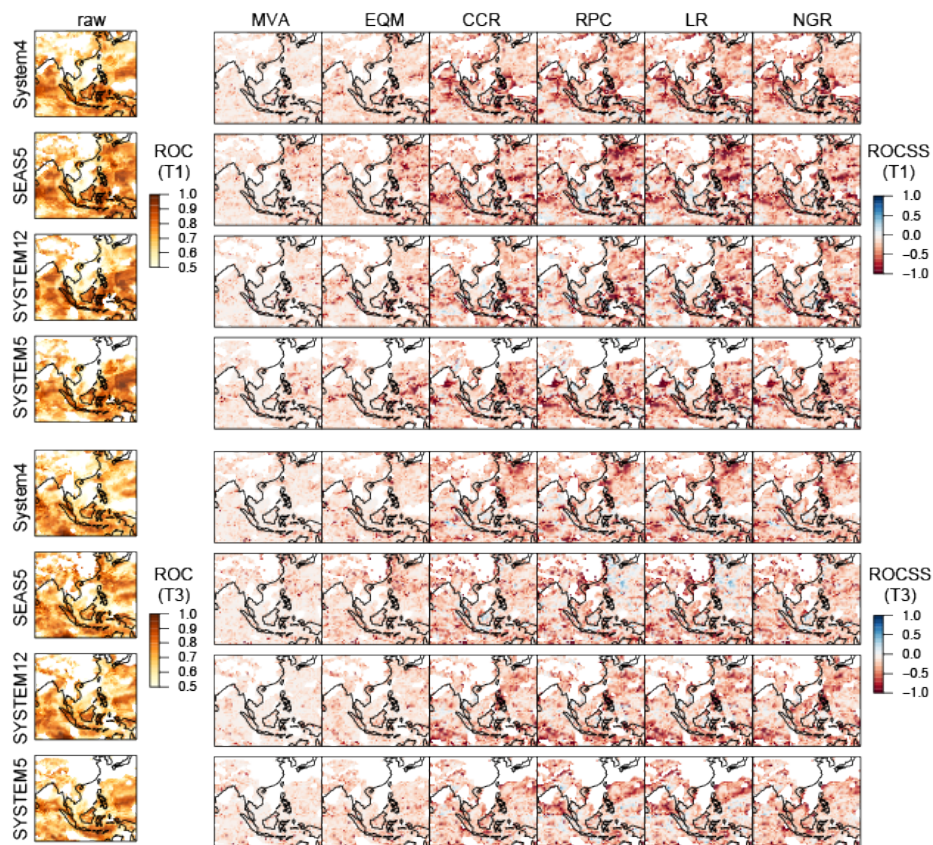
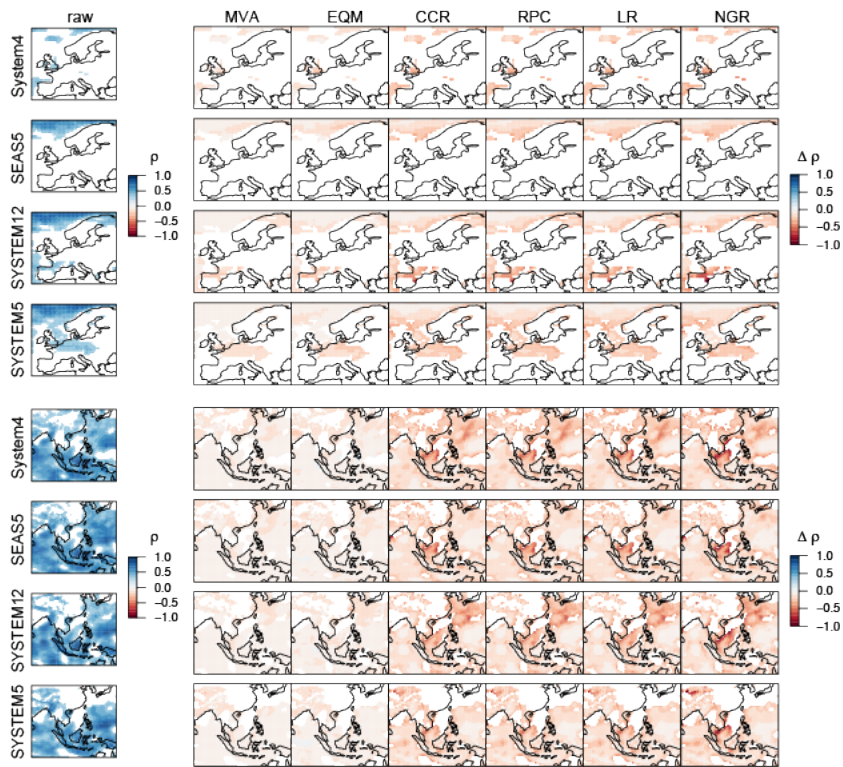


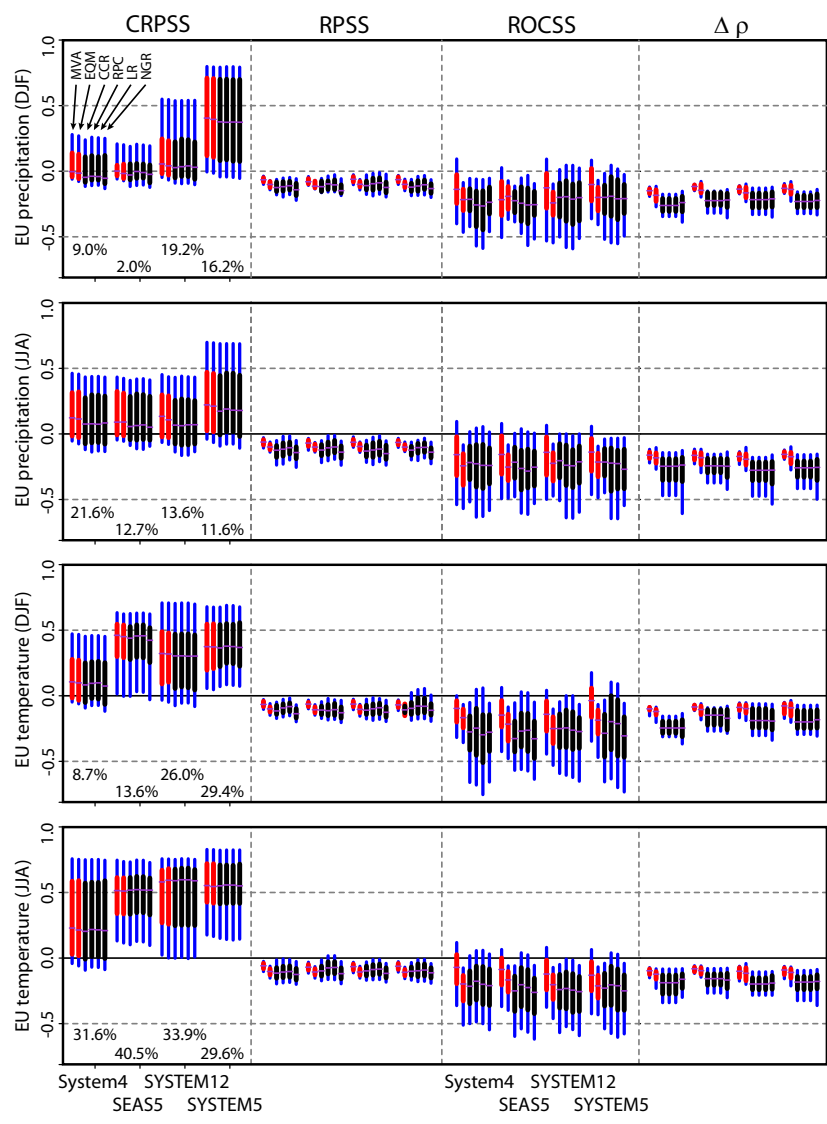
Fig. 5 Reliability and resolution components (top and bottom row, respectively) of the RPS for temperature over SA in DJF, as obtained from applying the BA and RC methods of Table 2 (in columns) to the System4 (12-member, 21-year version).



**Fig. 6** ROCSS for the cold (top) and warm (bottom) tercile categories of DJF temperature over SA, as obtained from applying the BA and RC methods of Table 1 (columns 2-7) to the four models of Table 1 (in rows). In all cases, the ROC obtained for the raw outputs (column 1) is considered as reference.



**Fig. 7** Column 1: Interannual Pearson correlation between ERA-Interim and the ensemble mean of the four available models (in rows) for DJF temperature, as given by the raw forecasts over EU (top) and SA (bottom). Columns 2-7: Difference (in correlation units) with respect to column 1, as obtained from the application of the BA and RC methods of Table 2.



**Fig. 8** Summary of the results obtained over EU, in terms of the different skill scores considered (CRPSS, RPSS, ROCSS and correlation differences; in columns). The two variables (precipitation and temperature) and seasons (DJF and JJA) analyzed are shown in different rows. In all cases, results for the four available models (System4, SEAS5, SYTEM12 and SYSTEM5) are displayed along the x-axis. For each model, the two (four) red (black) boxplots indicate the P25-75 range for each BA (RC) method, with blue corresponding to the P10-P90 range. The numbers in the first column correspond to the percentage of skillful gridboxes over which the methods were applied and tested (see Figure 2).



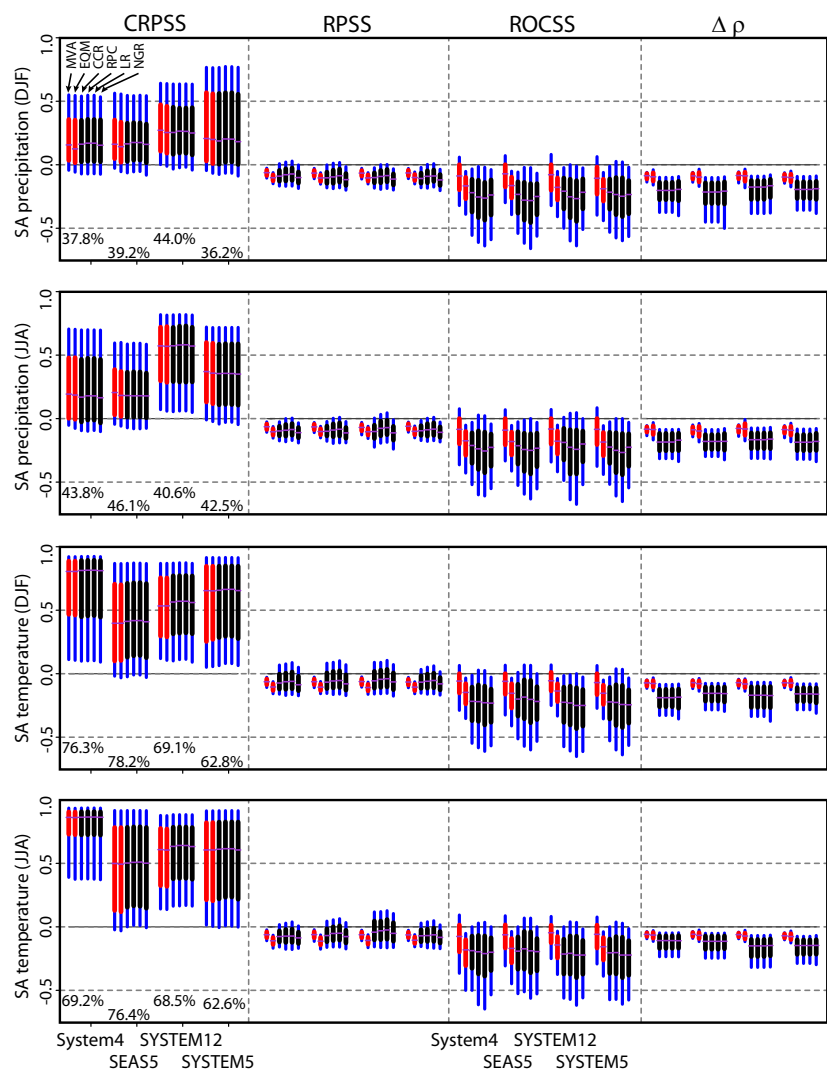
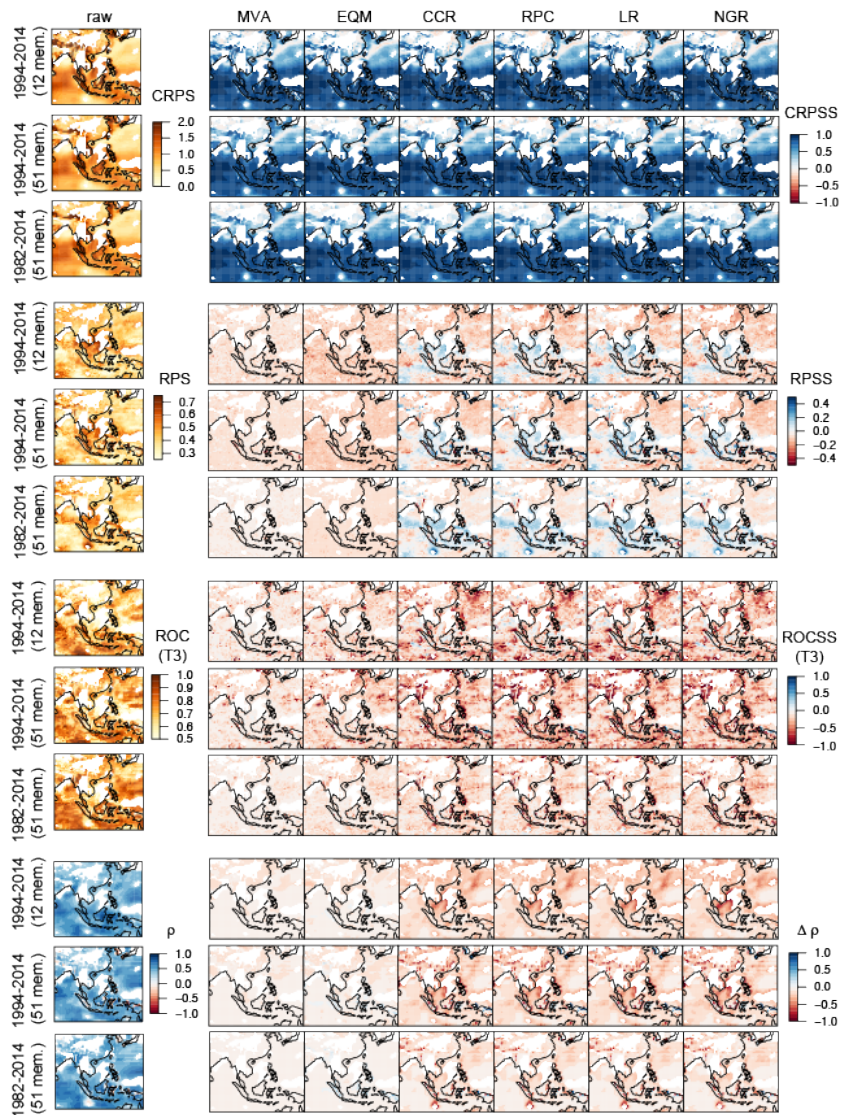
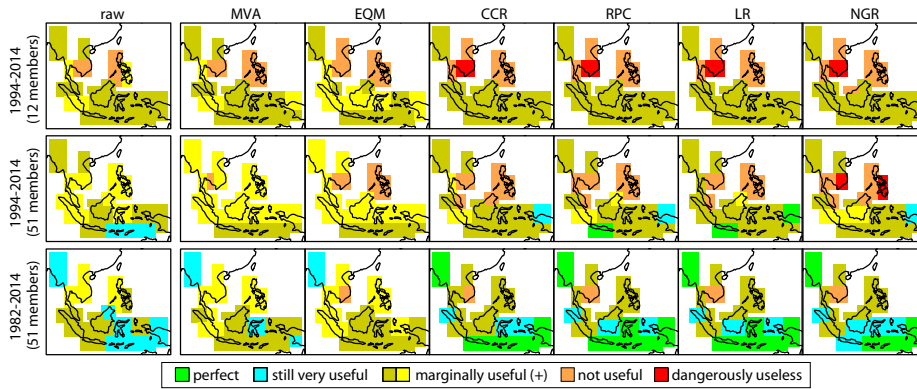


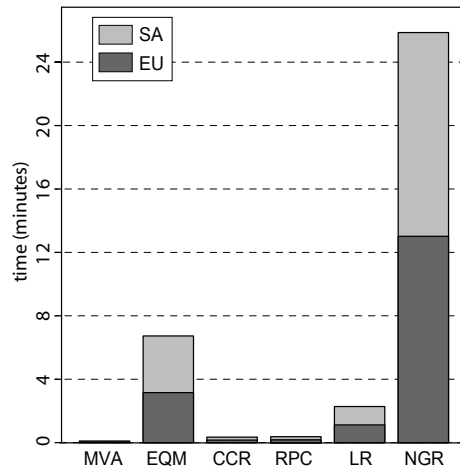
Fig. 9 As Figure 8, but for SA.



**Fig. 10** Results obtained for the CRPSS, the RPSS, the ROCSS (only for the warm tercile category) and the interannual Pearson correlation—in different panels from top to bottom—for temperature over SA in DJF, as obtained from applying the BA and RC methods of Table 2 to the System4. Within each panel, the top row corresponds to a 12-member ensemble for the period 1994-2014 (same as in Figures 3, 4, 6 and 7, displayed here again to facilitate comparison). The middle (bottom) row correspond to a 51-ensemble member for 1994-2014 (1982-2014).



**Fig. 11** Reliability categories, as obtained from applying the different BA and RC methods of Table 2 to correct DJF temperature from the System4 over SA. The top row corresponds to a 12-member ensemble for the period 1994-2014. Middle (bottom) row correspond to a 51-ensemble member for 1994-2014 (1982-2014).



**Fig. 12** Execution times—in a personal computer—for the different BA and RC methods of Table 2 for the illustrative case of temperature in DJF over EU (dark gray) and SA (light gray) for System4 (12-member and 21-year version), according to their implementation in the R-package *calibratoR*.